

Can Large Language Models Grasp Legal Theories? Enhance Legal Reasoning with Insights from Multi-Agent Collaboration

Anonymous ACL submission

Abstract

Large Language Models (LLMs) could struggle to fully understand legal theories and perform complex legal reasoning tasks. In this study, we introduce a challenging task (confusing charge prediction) to better evaluate LLMs' understanding of legal theories and reasoning capabilities. We also propose a novel framework: Multi-Agent framework for improving complex Legal Reasoning capability (MALR). MALR employs non-parametric learning, encouraging LLMs to automatically decompose complex legal tasks and mimic human learning process to extract insights from legal rules, helping LLMs better understand legal theories and enhance their legal reasoning abilities. Extensive experiments on multiple real-world datasets demonstrate that the proposed framework effectively addresses complex reasoning issues in practical scenarios, paving the way for more reliable applications in the legal domain.

1 Introduction

Large Language Models (LLMs) have shown remarkable generalization ability across diverse range of tasks and applications (Chowdhery et al., 2023; Touvron et al., 2023; OpenAI, 2023). But, current benchmarks may not adequately reflect the reasoning capabilities of LLMs (Valmeekam et al., 2024) and do not accurately reflect real-world situations (Huang and Chang, 2023). The validation of LLMs in more realistic and meaningful applications, such as legal reasoning, still requires extensive exploration.

In the legal domain, the core competency of legal professionals is to apply legal rules to facts and draw conclusions, as described by the IRAC (Issue, Rule, Application, Conclusion) framework. As shown in Figure 2, a legal professional can determine whether a case fact conforms to specific criminal charges based on legal rules. They critically assess a case against potential charges, focusing on

the key points of relevant legal rules, to accurately identify the appropriate charge and distinguish in-applicable charges. Legal rules, which manifest legal theories, determine the legal consequences of factual situations (MacCormick, 2005). Therefore, properly applying legal rules reflects the grasp of legal theories.

However, powerful LLMs may struggle to fully understand legal theories and perform basic legal reasoning tasks. Existing study (Dahl et al., 2024) has found that when LLMs are given criminal facts and legal rules, then asked whether cases constitute a certain charge, they tend to answer "yes," regardless of whether the charge is correct (golden charge) or a closely related one (confusing charge). Our empirical experiments also confirmed this issue. We sampled real-world criminal cases involving the charge of **Misappropriation of Public Fund**, inputting the criminal facts and legal rules into LLMs, and asked whether the case constituted the golden charge. Meanwhile, we created a control group where we input the same criminal facts and related legal rules, asking whether the case constituted a confusing charge (**Fund Misappropriation**). These two charges are very similar, with the key difference being *whether the defendant's subject position is that of a state functionary*. As shown in Figure 1, when performing legal reasoning, regardless of the prompt method or the version of GPT used, LLMs exhibit significant declines in performance when predicting confusing charges.

Generally, LLMs could face following challenges in legal reasoning: **Inconsistent reasoning**. Legal reasoning involves multi-step, compositional logic processes (Servantez et al., 2024). LLMs can be easily distracted by the interaction when generating reasoning steps (Shi et al., 2023) and may not be trustworthy by the tendency to give affirmative answers (Dahl et al., 2024). **Missing key details**. Legal rules and criminal facts are often described in complex natural language, making it

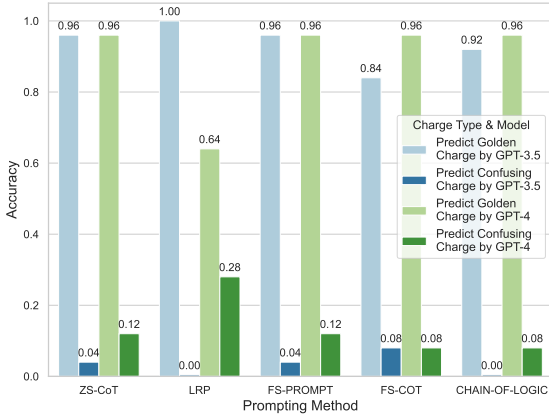


Figure 1: The performance of LLMs on predicting the golden (**Misappropriation of Public Fund**) or confusing charge (**Fund Misappropriation**) for the cases from CAIL-2018 datasets. The horizontal axis represents 5 advanced prompt methods to solve legal reasoning problems (detailed information is described in Section 5). In each method, GPT-3.5 and GPT-4 both exhibit a significant performance gap.

challenging for LLMs to fully understand and reason based on them. Consequently, they often overlook key information in the rules. **Lacking domain knowledge.** LLMs may hallucinate erroneous legal knowledge (Dahl et al., 2024) or encounter gaps in common-sense knowledge (Huang et al., 2023). Their overconfidence can obscure these shortcomings, making them difficult to identify (Ni et al., 2024).

To better evaluate LLMs’ understanding of legal theories and their reasoning capabilities, we introduce and construct a challenging task: confusing charge prediction (The detailed task definition is provided in Section 3). We also propose a novel framework: Multi-Agent framework for improving complex Legal Reasoning capability (MALR). First, an auto-planner breaks down complex legal rules into sub-tasks, allocating them to expert agents, reducing inconsistent reasoning in LLMs. Second, a non-parametric learning framework is proposed to draw adaptive rule-insights from trials and errors. To address the problem that LLMs may overlook crucial information in legal rules, we design a module that mimics human learning by gaining experience through reasoning trajectories and knowledge feedback, then learning insights through self-reflection. These insights supplement the rules, encouraging LLMs to focus on key factors from legal knowledge and fully understand the rules, while also guiding them to automatically

seek help when they feel uncertain. These designs effectively improve LLMs’ reasoning and critical-thinking skills.

Our contributions are threefold:

- We propose a multi-agent framework based on non-parametric learning, which encourages LLMs to automatically decompose complex legal tasks and extract insights from legal rules. Our framework assists LLMs in gaining a deeper understanding of legal rules and enhances their legal reasoning capabilities.
- We introduce a challenging task, predicting potentially confusing charges, to better evaluate LLMs’ understanding of legal theories and their reasoning capabilities.
- Extensive experiments are conducted on the multiple real-world datasets, demonstrating that the proposed framework can effectively addresses complex reasoning issues in real-world scenarios. Our work paves the way for more trustworthy application in legal domain¹.

2 Related Work

2.1 Legal AI and LLMs

Legal AI aims to improve legal tasks through AI techniques, particularly showing significant potential in alleviating the issue of “too many cases but too fewer legal experts” in the legal field (Katz et al., 2023; Dahl et al., 2024). One of the main challenges in legal domain is the training dataset can be considerably expensive and sparse (Sun et al., 2020), primarily comes in text, such as statutes, law articles and criminal cases. Under these circumstances, LLMs shows promising prospects in legal scenarios due to their powerful generalization capabilities in understanding and generating text. These applications include areas such as legal summarization (Deroy et al., 2023), legal document retrieval (Sun et al., 2024), legal question answering (Louis et al., 2024) and legal judgment prediction (Yu et al., 2022; Wu et al., 2023; Servantez et al., 2024).

2.2 Legal Reasoning and LLMs

Reasoning based on judicial rules and case fact descriptions is a fundamental ability of legal professionals, reflecting their understanding and application of legal theories (Servantez et al., 2024). Previous studies on Legal Judgment Prediction (LJP)

¹Code and data will be available after the double-blind review.

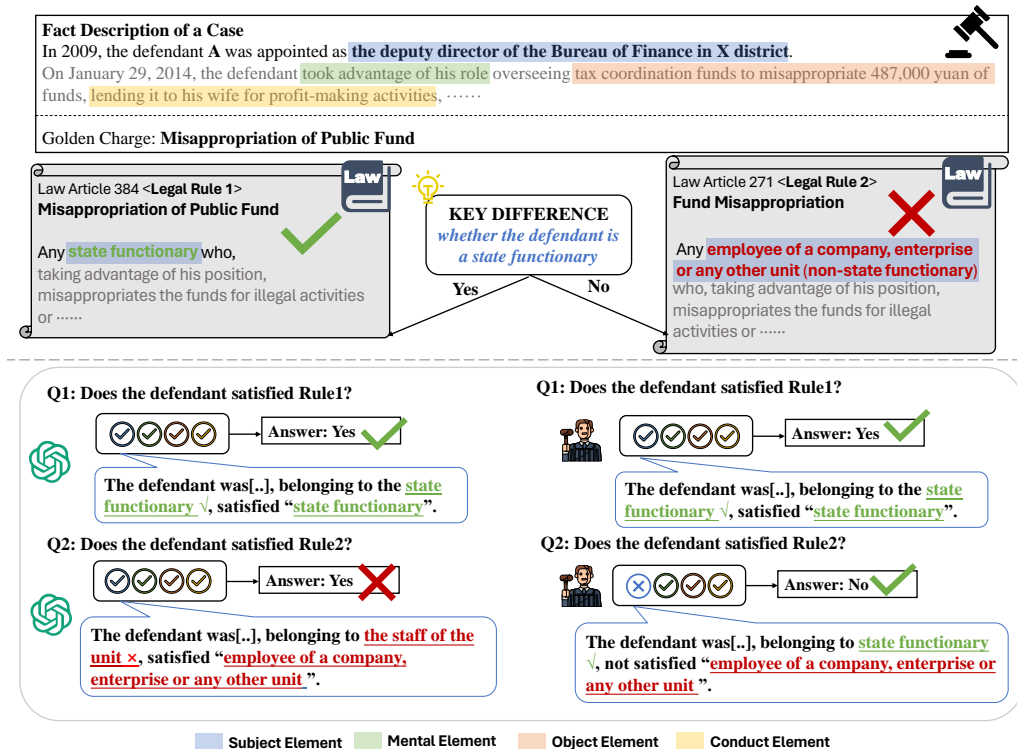


Figure 2: An example to demonstrate how a judge and an LLM apply legal rules to conclude whether a case satisfies a specific charge. This example outlines two confusing charges under Chinese criminal law: the **Crime of Fund Misappropriation** and the **Crime of Misappropriation of public fund**. The most significant difference between the two charges is **whether the defendant is a state functionary**. In the case description, the defendant is “*the Deputy Director of the Bureau of Finance in X district*”, a position that qualifies as a state functionary. Therefore, the judge can easily infer that the case falls under the Crime of Misappropriation of public funds, rather than Fund Misappropriation. However, the LLM fails to predict the confusing charge.

159 have primarily focused on automatically predicting
 160 case charges mainly from fact descriptions (Zhong
 161 et al., 2018; Chalkidis et al., 2019; Liu et al., 2022).
 162 Additionally, similar precedents can be retrieved as
 163 supplementary guides to improve performance (Wu
 164 et al., 2023). However, this approach can lead to
 165 inaccurate judgments due to overlooked potential
 166 differences in case details. To address the subtle
 167 differences between case details and legal rules,
 168 knowledge graphs have been introduced to solve
 169 confusing charges problems (Yue et al., 2021; Li
 170 et al., 2024). Despite these efforts, utilizing Four
 171 Elements Theory and innocent datasets, An et al.
 172 (2022) found that charge prediction models do not
 173 take legal theories into consideration. Instead, mod-
 174 els learn certain shortcuts for legal reasoning. Fur-
 175 thermore, Chain-of-Logic (Servantez et al., 2024)
 176 directly incorporates legal rules into prompts to
 177 elicit rule-based reasoning, achieving good perfor-
 178 mance on legal reasoning tasks involving three dis-
 179 tinct rules from the LegalBench benchmark. Simu-
 180 Court proposes a multi-agent framework to simu-

181 late the decision-making process of a judge (He
 182 et al., 2024).

183 Unlike existing works, we aim to evaluate and
 184 enhance the capacity of LLMs to reason based on
 185 legal theories, rather than treating legal rules as
 186 supplementary information.

3 Preliminary 187

188 We propose **Confusing Charge Prediction Task**
 189 to evaluate the LLMs’ ability to identify correct
 190 legal charges based on fact descriptions and legal
 191 rules, differentiating them from similar but incor-
 192 rect charges.

193 **Fact Description:** a concise description of a
 194 legal case, represented as a word sequence $\mathbf{f} =$
 195 $\{w_1, w_2, \dots, w_l\}$. **Legal Rule:** the definition of
 196 a specific criminal charge from law articles, also
 197 a word sequence $\mathbf{r}_c = \{w_1, w_2, \dots, w_n\}$, where c
 198 is the criminal charge. **Golden Charge:** The true
 199 crime label of a case. **Confusing Charge:** A charge
 200 similar to the golden charge but differing in one
 201 element (An et al., 2022).

To ensure LLMs’ trustworthiness in applying legal rules, we require them to confirm the golden charge as True and reject the confusing charge as False. The task can be formalized as:

$$y = \Gamma(f, r_{gc}) \wedge \neg\Gamma(f, r_{cc})$$

where gc refers to the golden charge, cc refers to the confusing charge, and Γ is the charge prediction model. y is True only if the fact description f satisfies the rule of golden charge r_{gc} and does not match the rule of confusing charge r_{cc} .

LLMs should correctly identify the golden charge and explain why the fact description doesn’t match the confusing charge, demonstrating understanding of legal theories.

4 The Proposed Framework

Figure 3 shows an overview of our proposed framework, which consists of four core components: Auto-Planner for Task Decompose, Role Assignment for Sub-task Agent, Adaptive Rule-Insights Training, and Reasoning with Rule-Insights.

4.1 Auto-Planner

A single LLM may exhibit inconsistencies when directly generating the whole reasoning process (Wang et al., 2024). Therefore, we designed an automatic planning module to decompose the task. Given a question q , a case fact description \mathbf{f} , and the corresponding legal rule \mathbf{r}_c about a criminal charge c , we guide an LLM as *auto-planner* to decompose the question into a sequence of sub-tasks based on the input of the fact and the rule:

$$[st_1, \dots, st_k] = LLM(q, \mathbf{r}_c, \mathbf{f}, p_{auto}) \quad (1)$$

where p_{auto} is the guideline prompt for LLMs to generate the sub-task set for the question q , and the st stands for the specific sub-task action and the k is the length of the sub-task set.

Given the resource-intensive nature of generating sub-tasks for every criminal cases, we design a more effective strategy. We first sample a smaller scale dataset, then generate the sub-task set for each sample. Subsequently, an LLM is used to identify semantically duplicate sub-task and compute the probability distribution for each sub-task. Based on this process, important sub-tasks with probability exceeding the threshold ζ are used to constitute the final sub-task set.

4.2 Assigning Roles to Sub-task Agent

Assigning roles can help the LLMs better perform complex task reasoning (Wang et al., 2023). Therefore, based on the sub-task set $[st_1, \dots, st_k]$, we employ k LLM-based agents to tackle each sub-task. Formally, we use the content of the sub-tasks to generate the appropriate prompts p_{st} and generate k agents to tackle each sub-task. Each agents will break down specific aspects of legal rule, check whether the fact description f conforms the corresponding sub-rule $r_{c_{st}}$ and generate the answer A_{st} . This process can be formulated as:

$$A_{st} = M_{st}(r_{c_{st}}, f, p_{st})$$

4.3 Adaptive Rule-Insights Training

As aforementioned, LLMs can be easily distracted by the irrelevant context (Shi et al., 2023) and tend to overlook the key information and important details within rules. Therefore, we aim to enable LLMs to automatically extract the most critical information for legal judgement directly from the rules. Previous research demonstrated that LLMs can autonomously learn from their own experiences (Shinn et al., 2024; Zhao et al., 2024). Inspired by the Kolb’s Experiential Learning Model (Kolb, 2014), we design the insights training module, as shown in Figure 3 (B), which consists of three core processes: experience gaining, insights drawing from errors and successes, and insights filtering. This module mimics human learning process and facilitates the LLMs to automatically learn rules, discovering and summarizing the most critical information in the rules.

Experience Gaining. A small training dataset with N charges is constructed, with each charge containing case samples and corresponding confusing charges. Based on the fact descriptions of a case, sub-task agents M_{st} will respectively generate sub-answers for both golden charge and confusing charge. These sub-answers will be synthesized into a final determination of whether the case satisfies the legal rule for the golden charge or confusing charge. On the l -th trial, ground truth is used as feedback. Successful trials are saved as successful experience, while failed trials trigger the Aspect-level Self-Reflector to identify incorrect sub-task agents M_{st}^e and generate reasons rs_e for the errors. In the next trial, the error reasons are used to improve sub-task agents’ predictions. Such approach of learning from trials and errors can be effective, as demonstrated in (Shinn et al., 2024).

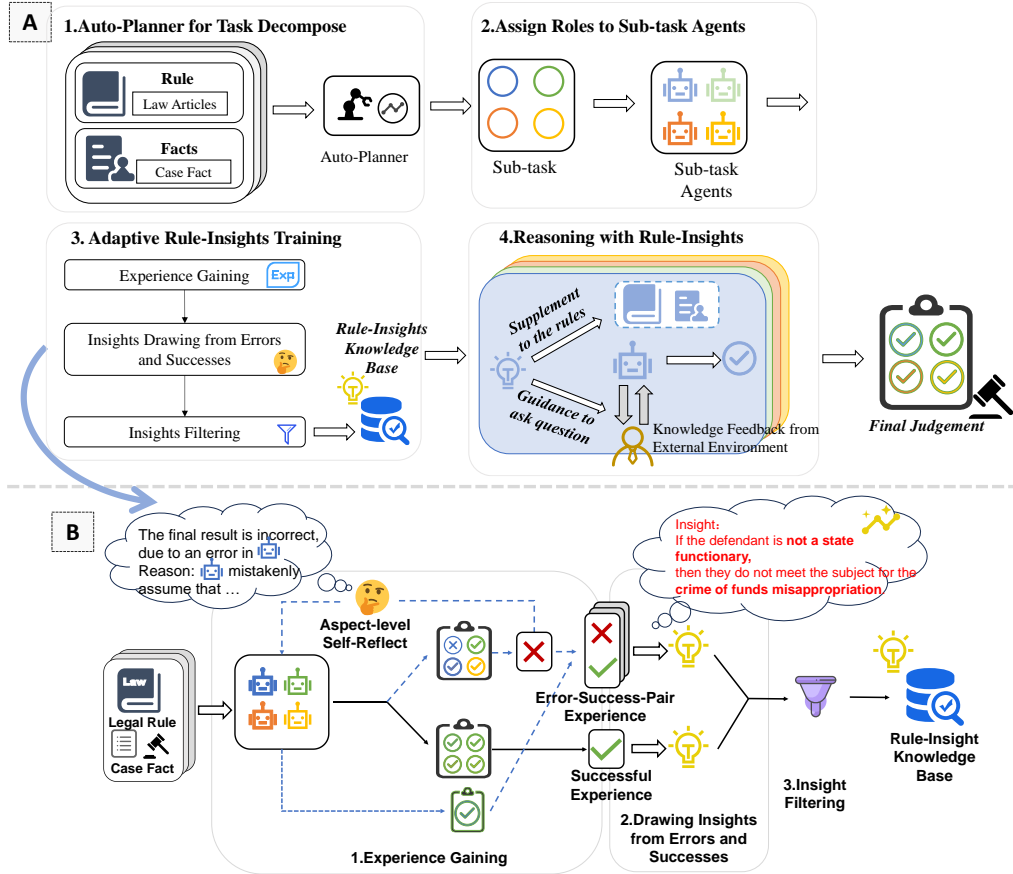


Figure 3: Our research framework in (A) and Adaptive Rule-Insights training process in (B).

This iterative process continues for a maximum of L rounds, and corrected trajectories are retained as error-success-pair experience. The algorithm procession is detailed in Alg 1.

Insights Drawing from Errors and Successes. We gain insights into rules by analyzing experience collections using different methods. For error-success pairs, we use a contrast-based approach, comparing incorrect and correct attempts. This enables the $M_{insight}$ to identify the most critical task-level judgments where errors occur, outputting insights in an if-then format. Successful experiences reveal common best practices (Zhao et al., 2024), so we provide the entire successful reasoning process to $M_{insight}$ to generate corresponding insights.

Insights Filtering. To address the potential for generating duplicate or incorrect insights when interpreting rules from the aforementioned process, we employ an LLM as an automatic checker, M_{filter} , to identify and remove redundant insights and filter out invalid ones. Ultimately, the retained insights are stored in the rule-insight knowledge base as

memory. The pseudo-code for insight drawing and filtering is presented in Algorithm 2.

4.4 Reasoning through Insights

The generated insights serve two purposes: (1) they supplement the rules as additional notes, and (2) they guide LLMs to inquire about uncertainties when facing knowledge gaps in specific sub-tasks.

For implementation, we retrieve relevant insights $inst$ from the knowledge base I for each question to improve reasoning. If the rule does not exist, the most similar rule from the knowledge base is selected based on semantic similarity, and a few-shot method is used to generate new insights. To address potential knowledge gaps in LLMs, our insights guide them to ask factuality questions. Based on the insights, we identify sub-tasks requiring fact-checking and use a few-shot method to prompt LLMs to ask key questions like “Is a <job position> a <state functionary>?” The generated question is then presented to a knowledgeable expert (a legal professional, a domain-specific LLM, or a search engine) to obtain knowledge feedback

Alg 1: Experience Gaining

Initialize: Self-Reflector, Sub-task Agent,
Evaluator: $M_{reflect}$, M_{st} , M_e
Number of charges N
Successful Experience $E_{success}$
Error-Success-Pair Experience E_{esp}
trajectory τ , Fact Description f
Golden Charge gc , Confusing Charge cc
while charge $n \leq N$ **do**
 Set $r_{gc}, r_{cc} \leftarrow gc_n, cc_n$
 Generate $\tau_{l,gc} = [A_{1,gc}, \dots, A_{k,gc}]$ using
 M_{st}, r_{gc}, f
 Generate $\tau_{l,cc} = [A_{1,cc}, \dots, A_{k,cc}]$ using
 M_{st}, r_{cc}, f
 while trial $l \leq L$ **do**
 Evaluate $(\tau_{l,gc}, \tau_{l,cc})$ using M_e
 if success **then**
 if $l = 1$ **then**
 Append $\tau_{l,gc}, \tau_{l,cc}$ to
 $E_{success}$
 break
 else
 Append $\tau_{l,gc}, \tau_{l,cc}$ to E_{esp}
 break
 end
 else
 Identify error M_{st}^e and Generate
 r_{se} using $M_{reflect}$
 $e \in \{gc, cc\}$
 Generate $A_{k,e}$ using
 M_{st}^e, r_e, f, r_{se}
 Updating $\tau_{l+1,e}$ using $A_{k,e}$
 end
 end
end

325 kg_{st} . Finally, we incorporate related insights in_{st}
326 and knowledge feedback kg_{st} into our ultimate rea-
327 soning process. As shown in Figure 3(A)(4), the
328 improved reasoning process for each sub-task agent
329 can be represented as:

$$A_{st} = M_{st}(r_{st}, f, in_{st}, kg_{st}, p_{st}) \quad (2)$$

331 All prompt templates for our MALR agents are
332 provided in Appendix A.

333 5 Experiment

334 5.1 Experimental Setting

335 **Dataset.** We evaluate LLMs' legal reasoning ca-
336 pability on three datasets: CAIL2018 (Xiao et al.,

Alg 2: Insight Drawing and Filtering

Initialize: Insight-Drawer, Insight-Filter:
 $M_{insight}$, M_{filter}
Successful Experience $E_{success}$
Error-Success-Pair Experience E_{esp}
Number of charges N
Number of sub-task k
Rule-Insight Knowledge Base I
while charge $n \leq N$ **do**
 Construct error-success pair of sub-task
 trial from E_{esp} :
 $\mathbf{P} = \{(A_{st_k}^{error}, A_{st_k}^{success}), \dots\}$
 for each p in \mathcal{P} **do**
 Drawing insight i using $M_{insight}(p)$
 Update i to $I[charge][st_k]$
 end
 for each exp in E_{esp} **do**
 Drawing insight i using
 $M_{insight}(exp)$
 Update i to $I[charge][st_k]$
 end
 Filter $I[charge]$ using M_{filter}
end

2018), CJO (Wu et al., 2023), and CAIL-I (An et al., 337
2022). CAIL2018 and CJO consist of real-world 338
cases with fact descriptions and golden charges. 339
We match corresponding confusing charges based 340
on the golden charges and randomly sample 400 341
cases from CAIL2018 and 100 from CJO for evalua- 342
tion. CAIL-I's testset contains 462 innocent cases 343
without crimes and the most similar charges to each 344
non-criminal fact. Further dataset information is 345
available in Appendix B. 346

The pairs of confusing charges are carefully 347
selected by a group of legal experts, including: 348
(1) Misappropriation of Public Fund (MP) v.s. 349
Fund Misappropriation (FM); (2) Bribery (BY) v.s. 350
Bribery of Non-State Officials (BN); (3) Kidnap- 351
ping (KD) v.s. Illegal Detention (ID); (4) Fraudu- 352
lently Obtaining Loans (FL) v.s. Loan Fraud (LF); 353
(5) Fund Misappropriation (FM) v.s. Official Em- 354
bezzlement (OE); (6) Fraud (FD) v.s. Loan Fraud 355
(LF); (7) Fraud (FD) v.s. Cheating and Bluffing 356
(CB); (8) Forging, Altering, Trading Official Docu- 357
ments, Certificates and Seals of State Organs (FO) 358
v.s. Forging the Seals of Companies, Enterprise, 359
Institution (FS). Key differences between each pair 360
are provided in Appendix B. 361

Implementation. We employ the publicly 362
available GPT-3.5-Turbo-0125 and GPT-4-0125- 363

364 preview models, with the temperature set to 0 for
 365 all text generation tasks. We sample two cases of
 366 each charge from the confusing pairs in CAIL-2018
 367 (totally 32 training cases) for auto-planner and in-
 368 sights training. The threshold ζ for the sub-task
 369 auto-planner is set to 0.8. Sentence-BERT (Thakur
 370 et al., 2021) and cosine similarity are used to com-
 371 pute semantic distances between rules and unseen
 372 legal rules, facilitating rule-insight inference test-
 373 ing in CJO and CAIL-I. During the insights training
 374 period, we limit the maximum number of trial at-
 375 tempts L to 2. For providing knowledge feedback,
 376 we employ Farui-200B², which can be replaced
 377 by other legal models or even legal experts in real-
 378 world scenarios. Additionally, we construct a legal
 379 rule knowledge base that includes Chinese Crimi-
 380 nal Law Articles and all charge definitions. All
 381 legal rules are retrieved from this knowledge base
 382 based on the charge name. The cost can be seen in
 383 Appendix B.

384 5.2 Baselines

385 **Zero-shot Setting:** (1) ZS-CoT (Kojima et al.,
 386 2022) uses “Let’s think step by step” to encour-
 387 ages LLMs to generate intermediate steps and im-
 388 prove reasoning. (2) Legal Reasoning Prompting
 389 (LRP) (Yu et al., 2022) is a zero-shot legal prompt-
 390 ing method that teaches LLMs to reason like a
 391 lawyer, following the “Approach, Issue, rule, appli-
 392 cation and conclusion” framework.

393 **Few-shot Setting:** (1) Few-Shot prompt-
 394 ing (Brown et al., 2020) is the standard prompting
 395 method includes only the sample and answer. We
 396 use a two-shot setting with one positive and one
 397 negative examples. (2) Few-Shot CoT (Wei et al.,
 398 2022) uses a few chain-of-thought demonstrations
 399 as exemplars to improve the ability of LLMs to
 400 perform complex reasoning. Again, we employ a
 401 two-shot setting with one positive and one negative
 402 examples. (3) Chain-of-Logic (Servantez et al.,
 403 2024) elicits rule-based reasoning by decomposing
 404 the rule into elements, answering each rule element
 405 separately, and finally using a logical expression
 406 to obtain the final answer. This approach is
 407 meticulously designed for legal reasoning tasks.

408 All prompt template can be seen in Appendix C.

409 5.3 Experiment Results

410 **Main Results:** From Table 2, we observe the fol-
 411 lowing findings. (1) LLMs fail to distinguish con-

412 fusing charges using simple but effective prompt
 413 methods such as CoT, and legal-specific prompt-
 414 ing approaches also fail to predict accurately. By
 415 examining the actual prediction results, we found
 416 that LLMs using these methods tend to respond
 417 with “yes.” (2) “MALR w/o insight”, which only
 418 decomposes the task into sub-tasks, outperforms
 419 all the baselines. This result indicates that de-
 420 composing the task into sub-tasks may mitigate
 421 LLMs’ biased tendencies. Notably, without any
 422 human intervention, our auto-planning strategy can
 423 decompose legal rules into four aspects: Subject
 424 (Sub), Mental (Men), Object (Obj) and Conduct
 425 (Con). This aligns with the Four Elements The-
 426 ory (An et al., 2022), which is widely recognized
 427 in the legal domain. (3) “MALR w/o ask” does
 428 not utilize external knowledge feedback but still
 429 achieves the second-best results, indicating that
 430 the learned insights did significantly enhance the
 431 LLM’s understanding of legal rules. (4) The com-
 432 plete MALR achieves the best performance on all
 433 datasets, demonstrating the effectiveness of pro-
 434 posed framework and the necessity of its core com-
 435 ponents. MALR achieves the best performance
 436 on nearly all confusing charge pairs (refer to Ap-
 437 pendix D). (5) Regarding the base models, GPT-3.5
 438 benefits more from our proposed MALR compared
 439 to GPT-4, achieving a more significant improve-
 440 ment over the baseline methods. This suggests that
 441 our framework has a stronger enhancing effect on
 442 LLMs with weaker foundational capabilities.

443 **Ablation Results:** Table 1 demonstrates the effec-
 444 tiveness of the components in adaptive rule-insights
 445 training module. (1) The results of “w/o $E_{success}$
 446 (without Successful Experience)”, “w/o E_{esp} (with-
 447 out Error-Success-Pair Experience)”, and “w/o
 448 $M_{filtering}$ (without Insight Filtering)” prove the
 449 significance of each designed component in the
 450 learning from the trial-and-error process. (2) The
 451 “directly generate” approach involves encouraging
 452 the LLM to generate insights directly based on the

Datasets Methods	CAIL2018		CJO	
	GPT-3.5	GPT-4	GPT-3.5	GPT-4
w/o insights	32.3	43.8	22.0	44.0
w/o $E_{success}$	38.8	50.0	29.0	48.0
w/o E_{esp}	46.0	48.8	33.0	48.0
w/o $M_{filtering}$	38.0	54.0	31.0	53.0
directly generate	32.0	43.3	35.0	38.0
complete MALR	40.8	56.8	39.0	55.0

Table 1: Ablation test for adaptive rule-insights training module.

²A legal-domain fine-tuned LLM based on Qwen (Bai et al., 2023), <https://tongyi.aliyun.com/farui>.

Methods	CAIL2018		CJO		CAIL-I	
	GPT-3.5	GPT-4	GPT-3.5	GPT-4	GPT-3.5	GPT-4
ZS-CoT (Kojima et al., 2022)	12.5	35.8	3.0	29.0	20.9	36.0
LRP (Yu et al., 2022)	9.8	37.8	1.0	37.0	22.3	49.6
FS-Prompt (Brown et al., 2020)	18.0	41.0	3.0	43.0	28.1	46.8
FS-CoT (Wei et al., 2022)	12.0	34.0	12.0	18.0	12.2	32.4
Chain-of-Logic (Servantez et al., 2024)	6.5	36.0	5.0	25.0	10.1	29.5
MALR w/o insight	32.3	43.8	22.0	44.0	45.3	53.2
MALR w/o ask	<u>37.3</u>	<u>53.3</u>	<u>31.0</u>	<u>53.0</u>	<u>51.1</u>	<u>55.4</u>
MALR (our)	40.8	56.8	39.0	55.0	56.1	57.6

Table 2: Main results on three legal datasets, the best is **bolded** and the second is underlined. The metric is accuracy. w/o insight refers to only decompose to sub-tasks, w/o ask refers to do not get any external knowledge feedback.

Fact Description of a Case In 2009, the defendant A was appointed as the deputy director of the Bureau of Finance in X district. On January 29, 2014, the defendant took advantage of his role overseeing tax coordination funds to misappropriate 487,000 yuan of funds, lending it to his wife for profit-making activities, [...]		
Golden Charge: Misappropriation of Public Fund ✓		Confusing Charge: Crime of Fund Misappropriation ✗
Legal Rule: Crime of Fund Misappropriation ✗ "Any employee of a company, enterprise or any other unit (non-state functionary) who, taking advantage of his position, misappropriates the funds for personal use or for loaning them to another person, or for profit-making activities or for illegal activities, [...]"		
[...] Issue: Does A commit the Crime of Fund Misappropriation? Rule: The essential elements for the crime include: 1. <u>The person is an employee of a company, enterprise, or other units;</u> 2. [...]; 3. [...]; 4. [...] Application: 1. A, <u>the deputy director of the Bureau of Finance in X district, is an employee of the unit;</u> 2. [...]; 3. [...]; 4. [...]. Conclusion: Based on the analysis above, A's actions meet the essential elements for the crime. Therefore, A is guilty of misappropriating funds. Final Answer: True .	[...] it is necessary to analyze from four elements: subject, mental, object, and conduct. Q1: Judgment on <u>whether the subject meets the criteria for the subject of the crime (staff members of companies, enterprises, or other units)?</u> A1: <u>X, as the deputy director of the Bureau of Finance in X district, is a staff member of a unit,</u> meeting the criteria for the subject of the crime of misappropriation of funds. (True) Q2: [...] ? A2: [...] (True) Q3: [...] ? A3: [...] (True) Q4: [...] ? A4: [...] (True) Based on the analysis, A constituted the crime. Final Answer: True	Judgment logic: According to the rule explanation, the specific subject of the crime of misappropriation of funds is the employee of a company, enterprise or any other unit, <u>and individuals with the status of state functionaries cannot be the subject of this crime.</u> In the case facts, A was appointed as the deputy director of the Bureau of Finance in X district. <u>Based on the nature of his position, A can be identified as a state functionary.</u> Therefore, as a state functionary, <u>A does not meet the specific subject requirements for the crime of Fund Misappropriation.</u> Answer: False [...] [...] [...]
LRP: True ✗	Chain-of-Logic: True ✗	MALR: False ✓

Figure 4: Case study for a given case. The **green** parts mean are the most critical information for distinguish the confusing charges, the **red** parts are contents that do not match the facts of the case.

legal rules without any training process. However, the performance drops in most situations, sometimes even worse than without using insights at all. A possible explanation is that directly generating insights may lead to the inclusion of unimportant information. We provide case examples with explanations comparing the directly generated insights with those obtained through our training process in Appendix E.

5.4 Case Study

Figure 4 presents an example of different methods used to predict confusing charges. As demonstrated in the case, our framework effectively focuses on the most critical aspects of the legal rules and makes a well-reasoned judgment. In contrast,

both LRP and Chain-of-Logic overlook the crucial information in the legal rules, resulting in their failure to accurately predict the confusing charge.

6 Conclusion

In the study, we introduce a challenging task to better evaluate LLMs' capability to comprehend legal theories. The proposed MALR framework can automatically decomposes complex legal tasks and extracts insights from legal rules, enhancing LLMs' legal reasoning abilities. Extensive experiments demonstrate MALR's effectiveness in equipping LLMs with a robust understanding of legal rules.

7 Ethical Considerations

The datasets we used for evaluation are all from public legal datasets, and information about the defendants has been anonymized. Thus, our study does not involve potential ethical concerns.

As Legal AI continues to evolve, our aim is to address the issue of "too many cases but too few legal experts." Moving forward, we hope to dedicate our efforts towards researching the interpretability of the LLMs, with the goal of building more reliable AI and supporting AI for good.

8 Limitations

Our work has two main limitations. First, even though we achieved great results, MALR did not predict correctly on all confusing charge pair cases. In the future, retrieval augmented generation could help our model perform better.

Second, our framework shows that LLMs can self-improve by summarize insights into the rules from trials and errors, which helps LLMs to better perform in complex legal reasoning tasks. Nevertheless, the potential for applying this approach in other fields such as medicine, finance, and scientific discovery remains unexplored. In the future, our framework could be applied in diverse domains.

References

Zhenwei An, Quzhe Huang, Cong Jiang, Yansong Feng, and Dongyan Zhao. 2022. [Do charge prediction models learn legal theory?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3757–3768, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

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*.

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

Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras. 2019. [Neural legal judgment prediction in English](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4317–4323, Florence, Italy. Association for Computational Linguistics.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul

Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.

Matthew Dahl, Varun Magesh, Mirac Suzgun, and Daniel E Ho. 2024. Large legal fictions: Profiling legal hallucinations in large language models. *arXiv preprint arXiv:2401.01301*.

Aniket Deroy, Kripabandhu Ghosh, and Saptarshi Ghosh. 2023. How ready are pre-trained abstractive models and llms for legal case judgement summarization? *arXiv preprint arXiv:2306.01248*.

Zhitao He, Pengfei Cao, Chenhao Wang, Zhuoran Jin, Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, and Jun Zhao. 2024. Simucourt: Building judicial decision-making agents with real-world judgement documents. *arXiv preprint arXiv:2403.02959*.

Jie Huang and Kevin Chen-Chuan Chang. 2023. [Towards reasoning in large language models: A survey](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1049–1065, Toronto, Canada. Association for Computational Linguistics.

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*.

Daniel Martin Katz, Dirk Hartung, Lauritz Gerlach, Abhik Jana, and Michael J Bommarito II. 2023. Natural language processing in the legal domain. *arXiv preprint arXiv:2302.12039*.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.

David A Kolb. 2014. *Experiential learning: Experience as the source of learning and development*. FT press.

Ang Li, Qiangchao Chen, Yiquan Wu, Xiang Zhou, Kun Kuang, Fei Wu, and Ming Cai. 2024. [From graph to word bag: Introducing domain knowledge to confusing charge prediction](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 7469–7479, Torino, Italia. ELRA and ICCL.

Dugang Liu, Weihao Du, Lei Li, Weike Pan, and Zhong Ming. 2022. Augmenting legal judgment prediction with contrastive case relations. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 2658–2667.

583	Antoine Louis, Gijs van Dijck, and Gerasimos Spanakis.	Karthik Valmeekam, Matthew Marquez, Alberto Olmo,	640
584	2024. Interpretable long-form legal question answer-	Sarath Sreedharan, and Subbarao Kambhampati.	641
585	ing with retrieval-augmented large language models.	2024. Planbench: An extensible benchmark for evalu-	642
586	In <i>Proceedings of the AAAI Conference on Artificial</i>	ating large language models on planning and reason-	643
587	<i>Intelligence</i> , volume 38, pages 22266–22275.	ing about change. <i>Advances in Neural Information</i>	644
		<i>Processing Systems</i> , 36.	645
588	Neil MacCormick. 2005. <i>Rhetoric and the rule of law:</i>		
589	<i>a theory of legal reasoning</i> . OUP Oxford.	Siyuan Wang, Zhongyu Wei, Yejin Choi, and Xiang Ren.	646
		2024. Can llms reason with rules? logic scaffolding	647
590	Shiyu Ni, Keping Bi, Jiafeng Guo, and Xueqi Cheng.	for stress-testing and improving llms. <i>arXiv preprint</i>	648
591	2024. When do llms need retrieval augmentation?	<i>arXiv:2402.11442</i> .	649
592	mitigating llms’ overconfidence helps retrieval aug-		
593	mentation. <i>arXiv preprint arXiv:2402.11457</i> .	Zekun Moore Wang, Zhongyuan Peng, Haoran Que,	650
		Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu,	651
594	OpenAI. 2023. Gpt-4 technical report. <i>arXiv preprint</i>	Hongcheng Guo, Ruitong Gan, Zehao Ni, Man	652
595	<i>arXiv:2303.08774</i> .	Zhang, et al. 2023. Rolellm: Benchmarking, elic-	653
		iting, and enhancing role-playing abilities of large	654
596	Sergio Servantez, Joe Barrow, Kristian Hammond, and	language models. <i>arXiv preprint arXiv:2310.00746</i> .	655
597	Rajiv Jain. 2024. Chain of logic: Rule-based reason-		
598	ing with large language models. <i>arXiv preprint</i>	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	656
599	<i>arXiv:2402.10400</i> .	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	657
		et al. 2022. Chain-of-thought prompting elicits reason-	658
600	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan	ing in large language models. <i>Advances in neural</i>	659
601	Scales, David Dohan, Ed H. Chi, Nathanael Schärli,	<i>information processing systems</i> , 35:24824–24837.	660
602	and Denny Zhou. 2023. Large language models can		
603	be easily distracted by irrelevant context . In <i>Proceed-</i>	Yiquan Wu, Siying Zhou, Yifei Liu, Weiming Lu, Xi-	661
604	<i>ings of the 40th International Conference on Machine</i>	aozhong Liu, Yating Zhang, Changlong Sun, Fei Wu,	662
605	<i>Learning</i> , volume 202 of <i>Proceedings of Machine</i>	and Kun Kuang. 2023. Precedent-enhanced legal	663
606	<i>Learning Research</i> , pages 31210–31227. PMLR.	judgment prediction with LLM and domain-model	664
		collaboration . In <i>Proceedings of the 2023 Conference</i>	665
607	Noah Shinn, Federico Cassano, Ashwin Gopinath,	<i>on Empirical Methods in Natural Language Process-</i>	666
608	Karthik Narasimhan, and Shunyu Yao. 2024. Re-	<i>ing</i> , pages 12060–12075, Singapore. Association for	667
609	flexion: Language agents with verbal reinforcement	Computational Linguistics.	668
610	learning. <i>Advances in Neural Information Process-</i>		
611	<i>ing Systems</i> , 36.	Chaojun Xiao, Haoxi Zhong, Zhipeng Guo, Cunchao Tu,	669
		Zhiyuan Liu, Maosong Sun, Yansong Feng, Xianpei	670
612	Changlong Sun, Yating Zhang, Xiaozhong Liu, and Fei	Han, Zhen Hu, Heng Wang, et al. 2018. Cail2018:	671
613	Wu. 2020. Legal intelligence: Algorithmic, data, and	A large-scale legal dataset for judgment prediction.	672
614	social challenges . In <i>Proceedings of the 43rd Inter-</i>	<i>arXiv preprint arXiv:1807.02478</i> .	673
615	<i>national ACM SIGIR Conference on Research and</i>		
616	<i>Development in Information Retrieval</i> , SIGIR ’20,	Fangyi Yu, Lee Quartey, and Frank Schilder. 2022. Le-	674
617	page 2464–2467, New York, NY, USA. Association	gal prompting: Teaching a language model to think	675
618	for Computing Machinery.	like a lawyer. <i>arXiv preprint arXiv:2212.01326</i> .	676
619	ZhongXiang Sun, Kepu Zhang, Weijie Yu, Haoyu Wang,	Linan Yue, Qi Liu, Binbin Jin, Han Wu, Kai Zhang,	677
620	and Jun Xu. 2024. Logic rules as explanations for	Yanqing An, Mingyue Cheng, Biao Yin, and Day-	678
621	legal case retrieval . In <i>Proceedings of the 2024 Joint</i>	ong Wu. 2021. Neurjudge: A circumstance-aware	679
622	<i>International Conference on Computational Linguis-</i>	neural framework for legal judgment prediction. In	680
623	<i>tics, Language Resources and Evaluation (LREC-</i>	<i>Proceedings of the 44th International ACM SIGIR</i>	681
624	<i>COLING 2024)</i> , pages 10747–10759, Torino, Italia.	<i>Conference on Research and Development in Infor-</i>	682
625	ELRA and ICCL.	<i>mation Retrieval</i> , pages 973–982.	683
626	Nandan Thakur, Nils Reimers, Johannes Daxenberger,	Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu	684
627	and Iryna Gurevych. 2021. Augmented SBERT: Data	Lin, Yong-Jin Liu, and Gao Huang. 2024. Expel:	685
628	augmentation method for improving bi-encoders for	Llm agents are experiential learners. In <i>Proceedings</i>	686
629	pairwise sentence scoring tasks . In <i>Proceedings of</i>	<i>of the AAAI Conference on Artificial Intelligence</i> ,	687
630	<i>the 2021 Conference of the North American Chapter</i>	volume 38, pages 19632–19642.	688
631	<i>of the Association for Computational Linguistics: Hu-</i>		
632	<i>man Language Technologies</i> , pages 296–310, Online.	Haoxi Zhong, Zhipeng Guo, Cunchao Tu, Chaojun Xiao,	689
633	Association for Computational Linguistics.	Zhiyuan Liu, and Maosong Sun. 2018. Legal judg-	690
		ment prediction via topological learning. In <i>Proceed-</i>	691
634	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	<i>ings of the 2018 conference on empirical methods in</i>	692
635	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	<i>natural language processing</i> , pages 3540–3549.	693
636	Baptiste Rozière, Naman Goyal, Eric Hambro,		
637	Faisal Azhar, et al. 2023. Llama: Open and effi-		
638	cient foundation language models. <i>arXiv preprint</i>		
639	<i>arXiv:2302.13971</i> .		

A Prompt Template for our MALR Agents

The prompt Templates for each agents can refer to Figure 5, 6, 7. We provide prompt templates in English; however, when applied in practice, these templates can be adapted to different languages by translating them into the corresponding language-specific prompts.

B Dataset and Experiments Information

CAIL2018 is a popular Chinese charge prediction datasets. It consists of real-world cases, each of which includes a fact description and the corresponding charges labels.

CJO is another Chinese legal dataset, same source from the CAIL2018, which is constructed to mitigate the potential data leakage.

CAIL-I contains 462 innocent cases that did not involve any crime. The dataset also has annotations for the criminal charge most similar to the non-criminal facts. The legal judgment prediction for an innocent case adheres to the presumption of innocence. It can evaluate whether LLMs can fully conform to legal rules for reasoning.

Key differences between each pair of confusing charges are provided in Figure 8. **Model Cost.** Statistically, the total token of our method is 1365 for each CAIL2018 example and the inference time per example is about 22s.

C Prompt Template for Baseline

The prompt templates for each baseline can refer to Figure 9, 10, 11. We provide prompt templates in English; however, when applied in practice, these templates can be adapted to different languages by translating them into the corresponding language-specific prompts.

D Specific Performance in the CAIL2018 dataset

Table 3 details the specific performance for each confusing-charge pair in the CAIL2018 dataset. The proposed MALR framework achieves the best performance on nearly all confusing charge pairs.

E More Cases for our insights

Figure 12 shows our training rule-insights can better learn the slight difference in the legal rules, which encourage the LLMs to better understand the legal rules.

Auto-Planner

You are currently in the task planning stage. Given a [Legal Rule Description] and related [Fact Descriptions of the case]. Please break it down into sub-tasks.

[Legal Rule Description]
{legal rule}

[Fact Descriptions of the case]
{fact description}

- Each sub-task action MUST have a unique ID, which is strictly increasing.
- Ensure the plan maximizes parallelizability.
- Never explain the sub-task actions with comments.

Sub-task Agent

You are a helpful legal profession. With a clear definition of the rule of {sub-task}.

Please determine whether {criminals} commit the crime of {charge_name} on the {sub-task} aspect, based on the [Legal Rule Description] and [Fact Descriptions of the case].

(Answer True if it constitutes the crime of {charge_name}, and answer False if it does not).

[Legal Rule Description]
{legal rule on sub-task}

Note: {rule-insights into the sub-task legal rule} //When training insights, this is Empty String//

[Fact Descriptions of the case]
{fact description}

[Knowledge feedback based on insight]

{Knowledge_feedback_by_external expert} //When training insights, this is Empty String//

Note:

Clarify the elements of {sub-task} and their corresponding relationship with the rules, clearly express your judgment logic, and provide a definite conclusion answer: True, False (answer True if it constitutes the {sub-task} of the crime of {charge_name}, answer False if it does not).

Answer format: [Judgment Logic] + [Answer]

Self-Reflector

You are an advanced legal agent who can analyze the incorrect answer and reasons through self-reflection.

By breaking down the task into following sub-tasks: {sub-task list}, sub-task experts reason that whether the defendant commits the crime of certain charge on the corresponding sub-aspect, based on the [Sub-task Legal Rule] and [Fact Descriptions of the case].

But sub-task experts incorrectly answered the question, please analyze where the judgment was mistaken based on the error trial, which could be one or more sub-tasks.

[Legal Rule Description]
{legal rule on sub-task}

[Fact Descriptions of the case]
{fact description}

[Incorrect Answer]
{initial_error_answers}

[ground truth]
{GROUND TRUTH FROM EXTERNAL FEEDBACK}

[Requirement]
[answer format]

Aspect1: <ONLY the option word of the four aspects; not a complete sentence!>

Reason1: <ONLY the reason why Aspect1 you conclude error results in Chinese>

...

Select the key error aspect, NOT all aspects are necessary to analyze.

Figure 5: Prompt Template for Auto-Planner, Sub-task Agent and Self-Reflector

Insight Drawer for error-success-pair experience

You are an advanced legal agent who can draw insight into the rule to improve by self-reflection. I will give your two attempts at answering a legal reasoning question based on a given the [Legal Rule Description] and [Fact Descriptions of the case].

There are one incorrect answer and one correct answer. Please generate one-sentence insight into the sub-task legal rule to highlight the most critical task-level judgment factor, NOT mention any specific information like defendant's name.

[Legal Rule Description]
{legal rule on error sub-task}

[Fact Descriptions of the case]
{fact description}

[Question]
Please determine whether {criminal} commit the crime of {charge_name} on the {sub-task} aspect, based on the [Legal Rule Description] and [Fact Descriptions of the case].

[Error Trial]
{error_trial}

[Success Trial]
{success_trial}

[Output]

Insight Drawer for successful experience

You are an advanced legal agent who can draw insight into the rule to improve by self-reflection. I will give your two attempts regarding the judgment of a case. The first is to answer whether the fact meets [Legal Rule of {golden_charge}], and the second is whether it meets [Legal Rule of {confuing_charge}].

Please generate one-sentence insight into the rule to highlight the most critical task-level judgment factor between the two charges. NOT mention any specific information like defendant's name.

[Fact Descriptions of the case]
{fact description}

[Legal Rule Description 1]
{Golden charge's legal rule}

[Question]
Please determine whether {criminal} commit the crime of {golden_charge}

[Answer]
{Successful Trial for all sub-tasks responses}

[Legal Rule Description 2]
{Confusing charge's legal rule}

[Question]
Please determine whether {criminal} commit the crime of {confusing_charge}

[Answer]
{Successful Trial for all sub-tasks responses}

[Output]

Figure 6: Prompt Template for Insight Drawer

Insight Filtering

You are an insight filtering who can filter the insights in the rule-insight knowledge base.

[Insights knowledge base]
{insight_from_knowledge_base} /JSON Format/

[Requirement]

1. Check the correctness for insights
2. Filter and remove duplicate insights
3. Don't change the original expression of any insights
4. Return the same json format as [Insights knowledge base]

Insight Inferencer

You are an expert at extracting the most critical information from rules, and you will be given some legal rules and the insights that have been extracted from them.

These insights can help judges make court decisions.

Please refer to the following rules and insights, and generate the corresponding insight within a new legal rule.

[Example 1]

Legal Rule:
{similar_rule}

Insight:
{similar_rule_insight}

[Your turn]
Legal Rule:
{new_charge_rule}
Insight:

Ask Key Question for Fact Checking

Please form a key question based on the [insight] and [case fact].

[Start of Examples]

[insight]

If the subject is a state functionary, it does not meet the subject criteria for the crime of fund misappropriation.

[case fact]

[...]The defendant, taking advantage of his position as a customer manager at the X of the Agricultural Bank of China XXX, misappropriated RMB 400,000 of the unit's funds under the name of loan customer XXX by forging materials required for the "second use of credit application" of a business loan in XXX name on January 6, 2015.

[...]

[Question]

Does Xiao Moujia qualify as the subject for the crime of fund misappropriation?

[Your response]

S1: Review of the subject for the crime of fund misappropriation: The defendant Xiao Moujia is a customer manager at X of the Agricultural Bank of China XXX.

S2: Relationship between the subject and the insight: If the subject is a state functionary, it does not meet the subject criteria for the crime of fund misappropriation.

S3: Therefore, the key question formed is: Is the customer manager at X of the Agricultural Bank of China XXX a state functionary?

[End of Examples]

[Your turn]

[insight] {insight}

[case fact] {fact}

[question] Does this case constitute the element of {charge_name}?

[Your response]

Figure 7: Prompt Template for Insight Filtering, Insight Inference and Ask Key Question for Fact Checking

Charge Name	Criminal Charge Full Name (Chinese Charge Name Translation)	Key Difference	label
MP	Misappropriation of Public Fund (挪用公款)	Whether the subject of the defendant is a state functionary	yes
FM	Fund Misappropriation (挪用资金)		no
BY	Bribery (受贿)	Whether the subject of the defendant is a state functionary	yes
BN	Bribery of Non-State Officials (非国家工作人员受贿)		no
KD	Kidnapping (绑架)	Whether the mental aspect is to extort property.	yes
ID	Illegal Detention (非法拘禁)		no
FL	Fraudulently Obtaining Loans (骗取贷款、票据承兑、金融票证)	Whether the mental aspect is aimed at illegal possession.	no
LF	Loan Fraud (贷款诈骗)		yes
FM	Fund Misappropriation (挪用资金)	Whether the mental aspect is aimed at illegal possession.	no
OE	Official Embezzlement (职务侵占)		yes
FD	Fraud (诈骗)	Whether the object is a property or loan.	property
LF	Loan Fraud (贷款诈骗)		loan
FD	Fraud (诈骗)	Whether the object is property or the credibility of a state authority.	property
CB	Cheating and Bluffing (招摇撞骗)		credibility
FO	Forging, Altering, Trading Official Documents, Certificates and Seals of State Organs (伪造、变造、买卖国家机关公文、证件、印章)	Whether the object (seal) belongs to a state institution.	yes
FS	Forging the Seals of Companies, Enterprise, Institution, or People's Organization (伪造公司、企业、事业单位、人民团体印章)		no

Figure 8: Key difference between each pair of confusing charge

Golden Charge	MP	FM	BY	BN	KD	ID	FL	LF	FM	OE	FD	LF	FD	CB	FO	FS
GPT-3.5																
ZS-CoT	4.0	0.0	12.0	8.0	32.0	4.0	0.0	0.0	0.0	16.0	36.0	0.0	24.0	8.0	56.0	0.0
LRP	0.0	0.0	0.0	0.0	32.0	0.0	4.0	0.0	0.0	16.0	20.0	4.0	20.0	4.0	48.0	8.0
FS-Prompt	0.0	0.0	0.0	0.0	40.0	8.0	0.0	0.0	8.0	8.0	76.0	4.0	72.0	4.0	68.0	0.0
FS-CoT	8.0	64.0	12.0	0.0	20.0	0.0	0.0	0.0	0.0	12.0	36.0	0.0	16.0	0.0	24.0	0.0
Chain-of-Logic	0.0	28.0	0.0	8.0	0.0	8.0	0.0	0.0	4.0	0.0	4.0	0.0	28.0	0.0	20.0	40.0
MALR (Our)	24.0	64.0	64.0	16.0	68.0	28.0	28.0	12.0	8.0	28.0	24.0	72.0	32.0	44.0	52.0	88.0
GPT-4																
ZS-CoT	12.0	52.0	68.0	12.0	24.0	0.0	0.0	0.0	4.0	40.0	96.0	4.0	96.0	8.0	76.0	80.0
LRP	20.0	76.0	60.0	32.0	16.0	44.0	8.0	0.0	28.0	24.0	80.0	8.0	80.0	16.0	56.0	60.0
FS-Prompt	12.0	56.0	84.0	32.0	20.0	0.0	16.0	60.0	88.0	20.0	56.0	32.0	92.0	20.0	40.0	28.0
FS-CoT	8.0	64.0	48.0	12.0	24.0	0.0	0.0	0.0	0.0	20.0	100.0	0.0	92.0	4.0	84.0	88.0
Chain-of-Logic	8.0	80.0	56.0	16.0	24.0	0.0	0.0	0.0	8.0	16.0	100.0	0.0	92.0	12.0	80.0	84.0
MALR (Our)	36.0	88.0	84.0	32.0	36.0	76.0	32.0	28.0	44.0	20.0	96.0	56.0	100.0	12.0	72.0	96.0

Table 3: Results on each criminal charge of confusing-charge pairs on CAIL2018 dataset.

ZS-CoT

You are a helpful legal profession.

Please determine whether {criminals} commit the crime of {charge_name} based on the [Legal Rule Description] and [Fact Descriptions of the case].

(Answer True if it constitutes the crime of {charge_name}, and answer False if it does not).

[Legal Rule Description]
{legal rule}

[Fact Descriptions of the case]
{fact description}

Let's think step by step.

ZS-LRP

You are a helpful legal profession.

Please determine whether {criminals} commit the crime of {charge_name} based on the [Legal Rule Description] and [Fact Descriptions of the case] through IRAC (Issue, Rule, Application, Conclusion) legal reasoning approach.

(Answer True if it constitutes the crime of {charge_name}, and answer False if it does not).

[Legal Rule Description]
{legal rule}

[Fact Descriptions of the case]
{fact description}

FS-Prompt

You are a helpful legal profession.

Please determine whether {criminals} commit the crime of {charge_name} based on the [Legal Rule Description] and [Fact Descriptions of the case].

(Answer True if it constitutes the crime of {charge_name}, and answer False if it does not).

[Legal Rule Description]
{legal rule}

Here are some demonstrations:

<Demonstration 1>

[Fact Descriptions of the case]

{fact description of the positive example}

[Question]: Whether {criminals_demo1} commit the crime of {charge_name}?

[Answer]:True

<Demonstration 2>

[Fact Descriptions of the case]

{fact description of the negative example}

[Question]: Whether {criminals_demo2} commit the crime of {charge_name}?

[Answer]:False

Now, it is your turn!

[Fact Descriptions of the case]

{fact description}

[Question]: Whether {criminals} commit the crime of {charge_name}?

[Answer]:

Figure 9: Prompt Template for baseline ZS-CoT, ZS-LRP and FS-Prompt

FS-CoT

You are a helpful legal profession.

Please determine whether {criminals} commit the crime of {charge_name} based on the [Legal Rule Description] and [Fact Descriptions of the case].

(Answer True if it constitutes the crime of {charge_name}, and answer False if it does not).

Here are some demonstrations:

<Demonstration 1>

[Legal Rule Description]

{legal rule}

[Fact Descriptions of the case]

{fact description of the positive example}

[Question]: Whether {criminals_demo1} commit the crime of {charge_name}?

[Judgment Logic]:

{chain_of_thought_for_demo1}

[Answer]: True

<Demonstration 2>

[Legal Rule Description]

{legal rule}

[Fact Descriptions of the case]

{fact description of the negative example}

[Question]: Whether {criminals_demo2} commit the crime of {charge_name}?

[Judgment Logic]:

{chain_of_thought_for_demo2}

[Answer]: False

Now, it is your turn!

[Legal Rule Description]

{legal rule}

[Fact Descriptions of the case]

{fact description}

[Question]: Whether {criminals} commit the crime of {charge_name}?

[Judgment Logic]:

Figure 10: Prompt Template for baseline FS-CoT

Chain-of-Logic

You are a helpful legal profession.

Please determine whether {criminals} commit the crime of {charge_name} based on the [Legal Rule Description] and [Fact Descriptions of the case].

(Answer True if it constitutes the crime of {charge_name}, and answer False if it does not).

Here are some demonstrations:

<Demonstration 1>

[Legal Rule Description]

{legal rule}

[Fact Descriptions of the case]

{fact description of the positive example}

[Question]: Whether {criminals_demo1} commit the crime of {charge_name}?

[Judgment Logic]:

Decompose the rule into elements:

The rule can be decomposed by (A) subject rule, (B) mental aspect rule, (C) object rule, (D) conduct aspect rule.

Logical Expression: (A and B and C and D)

Answer each rule element separately:

Q1: Does the defendant satisfy the subject rule (specific content in the subject rule of {charge_name})?

A1: The defendant is the xx, so satisfied the subject rule.(True)

Q2: Does the defendant satisfy the mental aspect rule (specific content in the mental aspect rule of {charge_name})?

A2: The defendant is the xx, so satisfied the mental aspect rule.(True)

Q3: Does the defendant satisfy the object rule (specific content in the object rule of {charge_name})?

A3: The defendant is the xx, so satisfied the object rule.(True)

Q4: Does the defendant satisfy the conduct aspect rule (specific content in the conduct aspect rule of {charge_name})?

A4: The defendant is the xx, so satisfied the conduct aspect rule.(True)

Logical expression with answer: (True and True and True and True) = True

So the defendant commits the crime of {charge_name}.

[Answer]The final answer is: True

<Demonstration 2>

[Legal Rule Description]

{legal rule}

[Fact Descriptions of the case]

{fact description of the positive example}

[Question]: Whether {criminals_demo2} commit the crime of {charge_name}?

[Judgment Logic]:

...

//The Judgment Logic format is similar to the Demonstration 1//

...

Logical expression with answer: (False and True and True and True) = False

So the defendant does not commits the crime of {charge_name}.

[Answer]The final answer is: False

Now, it is your turn!

[Legal Rule Description]

{legal rule}

[Fact Descriptions of the case]

{fact description}

[Question]: Whether {criminals} commit the crime of {charge_name}?

[Judgment Logic]:

Figure 11: Prompt Template for baseline Chain-of-Logic

Charge Name	Key Difference		Sub-task Legal Rule	Our Training Insights	Directly Generate Insights
Misappropriation of Public Fund	Whether the subject of the defendant is a state functionary	yes	Subject: The subject of this crime is a special subject, namely state functionaries.	If one is not a state functionary, then they do not meet the subject requirement for the crime of misappropriation of public funds.	pay attention to "Special subject, namely state functionaries"
Fund Misappropriation		no	Subject: The subject of this crime is a special subject, namely employees of companies, enterprises, or other units. Individuals with the status of state functionaries cannot be subjects of this crime.	If the individual is a state functionary, then they do not meet the subject requirement for the crime of funds misappropriation.	pay attention to "Employees of companies, enterprises, or other units"
Kidnapping	Whether the mental aspect is to extort property	yes	Mental aspect: This crime is subjectively constituted by direct intent, and has the purpose of extorting property or taking hostages.	If the action is not intended for the purpose of extorting property, then it does not meet the subjective requirement of the crime of kidnapping.	pay attention to 1: "Direct intent" 2: "The purpose of extorting property or taking hostages"
Illegal Detention		no	Mental aspect: The crime of illegal detention is subjectively characterized by intent and aimed at depriving another person of personal freedom.	If the main purpose of the perpetrator is to extort property, then it does not meet the subjective requirement of the crime of illegal detention.	pay attention to 1: "Intentionally" 2: "With the purpose of depriving another person of personal freedom" 3: "Negligence does not constitute the crime of illegal detention"

Figure 12: Case study for illustrating the effectiveness of our training insights.