

Delegating Legal Reasoning: An Agentic Approach to Judgment Prediction

¹Bhawna Kaushik, ²Yazdani hasan

1.bhawna.kaushik@niu.edu.in, Noida International University

2.yazhassid@gmail.com, Noida International University

Abstract

Legal judgment prediction (LJP) has become increasingly important in the legal field. In this paper, we identify that existing large language models (LLMs) have significant problems of insufficient reasoning due to a lack of legal knowledge. Therefore, we introduce GLARE, an agentic legal reasoning framework that dynamically acquires key legal knowledge by invoking different modules, thereby improving the breadth and depth of reasoning. Experiments conducted on the real-world dataset verify the effectiveness of our method. Furthermore, the reasoning chain generated during the analysis process can increase interpretability and provide the possibility for practical applications.

1 Introduction

Legal judgment prediction (LJP) is an important task in legal natural language processing (NLP), aims to make correct judgment predictions based on the case's fact description (Liu et al., 2023). The judgment predictions include law articles, charges, and terms of penalty (Xu et al., 2024). This task not only provides judgment references to lawyers and judges, as well as providing legal consulting services to the general public (Luo et al., 2017; Shulayeva et al., 2017; McGinnis and Pearce, 2013).

Recently, large reasoning models (LRMs) have made remarkable progress across in reasoning-

intensive tasks, including multi-hop question answering and strategic planning (Wan et al., 2024b; Choi et al., 2025). These models can perform multi-step reasoning that mimics human thinking (Fue et al., 2022). Intuitively, LJP appears to be an ideal fit for such models. Legal decision-making often involves comparing multiple candidate charges, evaluating whether each satisfies the legal criteria, and narrowing down to the most appropriate one based on the case facts. As a result, it is natural to expect that strong reasoning models would lead to major improvements in LJP.

However, existing reasoning models fail to deliver the expected breakthroughs in LJP. In practice, they tend to predict the most likely charges without comparing them to similar alternatives, and their reasoning chains are often short and lacking in meaningful intermediate steps. These issues become especially clear in cases involving rare or confusing charges, where accurate judgment depends on subtle distinctions and careful reasoning. Although models may produce step-by-step outputs in such scenarios, the reasoning often stays at a surface level, focusing on pattern matching rather than legal principles.

We argue that the main reason for the limited performance of reasoning models in legal judgment tasks is not a lack of reasoning ability, but a lack of the specialized knowledge that legal reasoning depends on (Yuan et al., 2024). Effective legal analysis requires long-tail legal knowledge, such as determining the applicability of specific statutes. In some cases, this knowledge is even absent from official legal texts. When such information is missing, models struggle to produce complete and trustworthy reasoning chains as shown in Figure 1. These observations highlight the need for domain-specific

knowledge augmentation mechanisms that can dynamically supply essential information during the reasoning process.

To address the knowledge gaps in legal reasoning, we propose GLARE (*AGentic Legal Reasoning Framework*), a modular system that

International Research Journal of Multidisciplinary Sciences
VOL-1 ISSUE-7 July 2025 PP:1-10

enables language models to dynamically acquire key legal knowledge to improve the breadth and depth of reasoning. First, the **Charge Expansion Module** (CEM) expands a diverse set of confusing charges by leveraging multiple signals, such as legal structure and historical co-occurrence. This

ISSN:AWAITED

helps the model compare a wide range of candi-

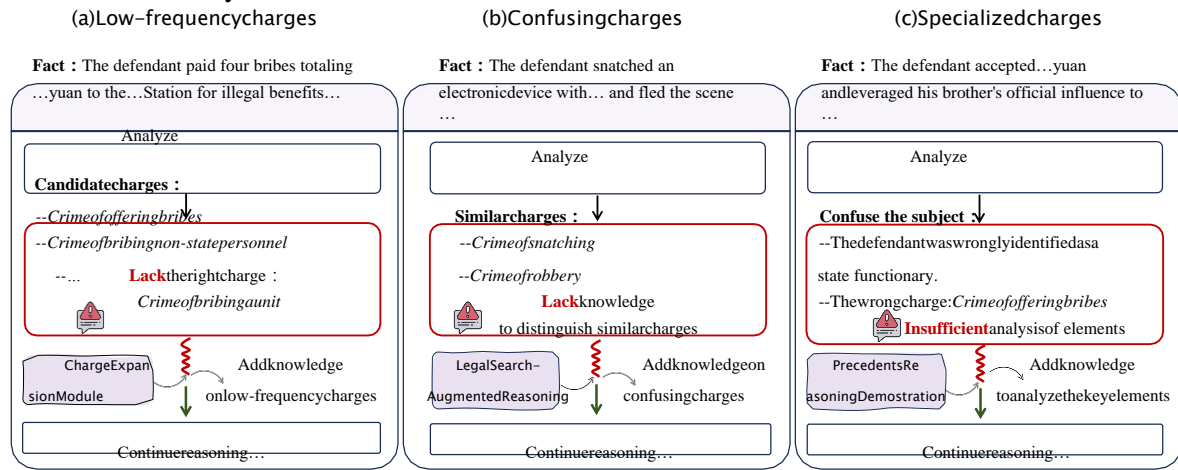


Figure1: Lack of knowledge in three aspects: (a) Lack of knowledge of low-frequency charges. (b) Lack of knowledge of confusing charges. (c) Lack of knowledge to analyze the key elements of the charges with strong professionalism.

dates and avoid premature conclusions. Second, the **Precedents Reasoning Demonstration (PRD)** module is built on reasoning paths that are constructed offline from real legal cases. During inference, the model retrieves the most relevant precedents through semantic search and learns from their reasoning chains via in-context learning. Finally, the **Legal Search-Augmented Reasoning (LSAR)** module allows the model to detect knowledge gaps and retrieve supporting legal information when needed. We guide the model to focus its search on differences between similar charges and details of how specific laws apply, rather than general case facts. Retrieved content is structured and injected into the reasoning process to support more accurate conclusions. By integrating essential legal knowledge, the model achieves more trustworthy and transparent judgment prediction.

Following prior work in legal judgment prediction, we conduct experiments on two publicly available real-world legal datasets. Experimental results show that our method consistently outperforms a range of strong baselines. Notably, it achieves substantial improvements on challenging cases involving confusing and difficult charges, where long-tail legal knowledge is crucial. These gains stem from our approach's ability to effectively enrich and incorporate relevant legal knowledge.

In summary, our contributions are as follows:

(1) We introduce **GLARE**, an agentic framework for legal judgment prediction that enhances reasoning by dynamically integrating legal knowledge throughout the decision-making process.

(2) We design three complementary modules to enrich the model's reasoning process by expanding candidate charges, leveraging real-world precedents, and injecting retrieved legal knowledge.

(3) Extensive experiments on two real-world datasets show that **GLARE** significantly outperforms strong baselines, with especially notable gains on cases requiring crucial legal knowledge.

2 Related Work

Legal judgment prediction Legal judgment prediction has experienced significant development and become an increasingly crucial NLP task. Earlier research (Segal, 1984) relied on artificially designed features to capture information from legal texts. Sulea et al., 2017 applied traditional machine learning methods to predict the legal judgment. Recent advances in deep learning (Xu et al., 2020; Zhang and Dou, 2023) have motivated researchers to leverage neural networks for automated text representation learning. Recently, LLMs have further promoted the progress of LJP (Denget al., 2024a), and several studies (Wu et al., 2023; Peng and Chen, 2024) employ Retrieval-Augmented Generation (RAG) technology (Zhao et al., 2024) to enhance LLMs by incorporating external legal knowledge. However, existing LLM-based methods struggle to utilize comprehensive legal knowledge (Fei et al., 2023) and refer to the way of precedent reasoning to analyze cases. In this context, we make full use of external knowledge and precedents.

Reasoning skills in language models Recent work has improved LLMs' reasoning through better prompting techniques (Sahoo et al., 2024). Wei et al. (2022) showed that chain-of-thought prompting

can explicitly guide LLMs to reason step by step. In the legal domain specifically, LoT ([Jiang and Yang, 2023](#)) proposed legal syllogism reasoning to improve performance on LJP

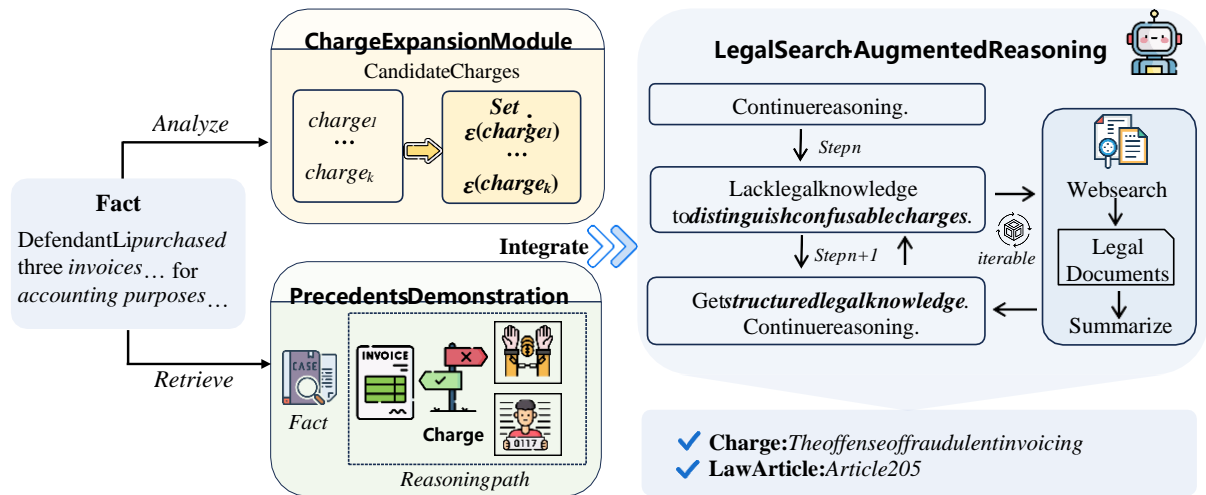


Figure2: Overview of our agentic legal reasoning framework. LLMs can utilize three external modules to acquire knowledge: ChargeExpandModule expands a diverse set of charges, precedents retrieved from an offline built database can provide in-context learning, Legal Search-Augmented Reasoning allows the model to detect knowledge gaps and retrieve supporting legal information.

task. ADAPT (Deng et al., 2024b) further established a comprehensive workflow for LJP that enables discriminative reasoning in LLMs. However, these approaches primarily rely on the LLMs' intrinsic capabilities, which inherently constrain the reasoning breadth and the depth of analysis (Zhang, 2024; Ke et al., 2025). Therefore, we propose an agentic legal reasoning framework to dynamically acquire key legal knowledge to improve the breadth and depth of reasoning.

3 Methodology

Preliminaries

We first formally define legal judgment prediction. Given a case fact description f , the model will analyze and predict the final judgment results including the relevant law articles, the convicted charges and the term of imprisonment for the defendant. Following previous works (Shui et al., 2023; Wei et al., 2025), we exclude the task of sentencing prediction from our scope as its subjective nature brings challenges that are not well aligned with the current capabilities of large language models. In this work, we treat large language models as agentic legal reasoners that can dynamically acquire and incorporate external legal knowledge to enhance their analysis. Rather than relying solely on parametric knowledge, our approach equips the model with access to external modules, enabling it to enrich its reasoning with case-specific legal

knowledge. ChargeExpandModule expands a diverse set of charges, precedents retrieved from an offline built database can provide in-context learning, Legal Search-Augmented Reasoning allows the model to detect knowledge gaps and retrieve supporting legal information. Given a case fact description f and a set of external modules M , the model performs step-by-step analysis to construct a coherent reasoning chain R and arrive at a final judgment prediction p . We formalize this process as a mapping: $(f, M) \rightarrow (R, p)$.

Agentic Legal Reasoning Framework

We propose GLARE, an agentic legal reasoning framework that autonomously invokes external modules to support comprehensive and informed judgment prediction. As shown in Figure 2, GLARE follows a structured three-stage reasoning pipeline:

1. **Charge Expansion:** The model begins by analyzing the case facts and generating preliminary candidate charges. To prevent premature narrowing of the decision space, it triggers the Charge Expansion Module to supplement the initial candidates with legally similar charges.
2. **Precedent-Enhanced Reasoning:** The model retrieves relevant precedents from an offline-constructed database that includes fact descriptions and synthesized reasoning chains. These reasoning chains were constructed in advance to illustrate the key distinctions between confusing charges. These precedents serve as case-specific reasoning demonstrations, helping the model better understand how similar legal criteria apply and guiding it through more precise reasoning via in-context learning.
3. **Iterative Search-augmented Reasoning:** As the model reasons through each candidate charge,

it dynamically identifies knowledge gaps such as missing legal definitions and charge-specific thresholds. Rather than treating retrieval as a one-time step, the model interleaves reasoning and retrieval in a loop. Retrieved results are injected back into the reasoning context, enabling the model to refine its current analysis. This iterative process continues until the model has collected sufficient knowledge to complete its reasoning and reach a final judgment.

The three modules collaboratively supplement legal knowledge and extend the legal reasoning chain. Next, we will introduce these three modules in detail.

Charge Expansion Module

To enable charge comparison and avoid premature conclusions, we expand each candidate charge by retrieving related charges. The expansion is based on two complementary perspectives: legal structure and historical co-occurrence.

Legal Structure-based Expansion. The Criminal Law is organized into chapters, each representing a specific legal interest or domain. Charges within the same chapter typically differ in subtle legal criteria, while charges across different chapters may involve similar actions or consequences but fall under distinct legal categories. To capture both fine-grained intra-domain distinctions and cross-domain conceptual similarities, we retrieve related charges from both within the same chapter and across different chapters.

Specifically, for a given charge c , we use the pretrained denser retriever BGE (Xiao et al., 2024) to find the top- k most similar charges from (a) the same chapter and (b) other chapters:

$E_1(c) = \text{topk}_{\text{same}}(c) \cup \text{topk}_{\text{diff}}(c), (1)$ where $\text{topk}_{\text{same}}(c)$ and $\text{topk}_{\text{diff}}(c)$ represent the most similar charges from the same and different chapters, respectively. This dual-source expansion helps the model compare similar alternatives, reducing the risk of overlooking relevant charges.

History-based Expansion. Certain charges tend to appear together in real-world cases, reflecting practical legal dependencies or common joint indictments. We leverage the MultiLJP (Lyu et al.,

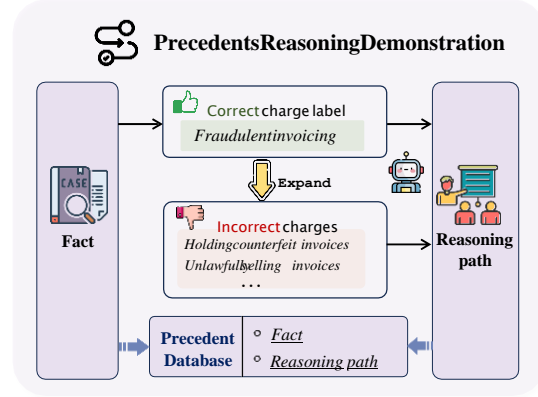


Figure 3: The module of Precedents Reasoning Demonstration: LLM analyzes the reasons for the selection or exclusion of each charge based on facts, thereby generating the reasoning path of precedents.

2023) dataset, where each case may involve multiple defendants and multiple charges. By analyzing these cases, we construct a co-occurrence dictionary that records how frequently each pair of charges appears together. For a given charge c , we select the top- k most frequently co-occurring charges as the expansion set $E_2(c)$.

Final Expansion Set. Given an initial set of candidate charges $\{c_1, c_2, \dots\}$ predicted by the language model, we apply the two strategies above to expand each charge:

$$E(c_i) = E_1(c_i) \cup E_2(c_i) \quad (2)$$

Precedents Reasoning Demonstration

Previous precedent-based approaches (Wu et al., 2023; Chen and Zhang, 2023; Santos et al., 2024) typically retrieve the fact description and final judgment of prior cases, then insert them directly into the prompt. However, such methods offer little insight into the reasoning process behind those decisions. As a result, they tend to rely on shallow fact matching rather than learning how to distinguish between legally similar charges.

To address this issue, we construct reasoning-augmented precedents that make the decision logic explicit. As shown in Figure 3, we first expand the original charge c into a set of similar charges C . Given the case fact f , the correct charge c , and the set of alternatives C , we prompt LLM to generate a reasoning path that explains why c is ap-

proprateandwhytheothercandidatesin $C\{c\}$ shouldbeexcluded¹. Thisreasoningisgenerated **offline**andstoredtogetherwiththecasefacts.

LegalSearch-AugmentedReasoning

While recent retrieval-augmented generation (RAG)approaches(Wuetal.,2023;PengandChen, 2024;Fengetal.,2024)enhancelegalmodelsbyretrieving precedents, statutes, and charge definitions, they remain limited in key aspects. Specifically, they often fail to resolve fine-grained distinctions between similar charges or providedetailedrulesto determine facts. Moreover, these methods rely on static retrieval from fixed knowledge bases, making them inflexible and unable to accommodate evolving judicial practices.

To address these limitations, we introduce a **dynamic and iterative legal search-augmented reasoning mechanism**. Rather than passively injecting generic legal reasoning and generating targeted queries. These queries focus on *subtledifferences between candidate charges* and *fact-specific questions*. We exclusively source authoritative legal interpretations from official channels, thereby minimizing noise. The system retrieves relevant legal texts from the web in real time, enabling up-to-date and context-related augmentation.

We further ground the model’s reasoning in a **sylogistic structure**: the retrieved legal context serves as the major premise, the case fact as the minor premise, and the conclusion is derived through logical alignment (Jiang and Yang, 2023; He et al., 2025). This structure helps the model remain grounded in factual evidence and reduce hallucinations. The overall reasoning process is formalized as an iterative function:

$$R_t = f_\theta(R_{<t}, q_t, d_t, f), \quad (3)$$

where R_t denotes the current reasoning state, $R_{<t}$ are the historical reasoning paths, q_t and d_t are the query and corresponding retrieved documents of this step, and f is the case fact.

This design enables the model to incrementally construct a legally grounded reasoning chain, adapting

content, our method allows the LLM to actively identify knowledge gaps during needs, our framework offers greater flexibility to real-world legal dynamics.

4 Experiments

Datasets and Evaluation

We conducted experiments in both single-defendant and multi-defendant scenarios to verify the effectiveness of our method in practical applications. For the single-defendant case, we use the CAIL2018 dataset (Xiao et al., 2018). For the multi-defendant case, we adopt the CMDL dataset (Huang et al., 2024). We uniformly sampled across all charges to construct a balanced test set. The details are shown in Table 1. For the PRD module, We employ the training set from both dataset as our precedent database. For evaluation metrics, we adopt the same measures used in prior work: Accuracy (Acc.), Macro Precision (Ma-P), Macro Recall (Ma-R), and Macro F1 (Ma-F).

tively integrating external knowledge as needed. By decoupling retrieval from static knowledge bases and aligning it with the model’s evolving

Dataset	CAIL2018	CMDL
#Train cases	100,531	63,032
#Test cases	1,000	834
#Charges	190	164
#Articles	175	147
#Average criminal per case	1	3.79
Average length per case	409.6	1124.94

Table 1: Statistics of dataset.

Baselines

We compare our method against two categories of baseline approaches:

Classification Methods: These methods take legal judgment prediction as a classification task, relying on supervised learning with labeled datasets.

TopJudge (Zhong et al., 2018) employs a graph structure to model the topological dependency among the three subtasks: charge prediction, law article prediction, and sentence term prediction.

¹ We provide the detailed prompt and examples for synthesizing reasoning paths in Appendix C

NeurJudge (Yue et al., 2021) integrates a legal knowledge graph into the neural architecture, capturing explicit relationships among legal entities and improving reasoning over structured legal knowledge. **BERT** (Devlin et al., 2019), a standard pre-trained transformer model, is adapted to legal texts via supervised training. It serves as a strong baseline for judgment prediction tasks. **Lawformer** (Xiao et al., 2021) is built upon Longformer (Beltagy et al., 2020) and further pretrained

(2) *Search-o1* (Liet al., 2025) dynamically retrieves external knowledge when it encounters uncertain or ambiguous knowledge in the general domain. We use reasoning model QwQ-32B (Team, 2025) in this setting.

Experiment Settings

In our experiments, we adopt Qwen2.5-32B (Yang et al., 2024a) and QwQ-32B (Team, 2025) as the

Methods	Charge				LawArticle			
	Acc.	Ma-P	Ma-R	Ma-F	Acc.	Ma-P	Ma-R	Ma-F
Classification Methods								
TopJudge	52.1	50.9	45.7	43.5	52.8	47.7	43.8	41.2
LADAN	76.7	73.4	71.0	69.5	77.5	71.0	69.2	67.5
NeurJudge	74.7	77.7	71.5	71.5	77.4	80.7	74.6	74.3
BERT	85.8	83.4	86.6	83.3	85.8	80.4	82.8	79.9
Lawformer	71.3	58.2	62.7	57.8	72.9	58.1	61.4	56.9
Direct Reasoning								
Qwen2.5-32B	74.5	75.3	69.3	69.1	77.1	73.3	66.6	67.1
QwQ-32B	82.5	86.9	80.5	80.9	84.0	83.1	76.1	77.0
Qwen2.5-72B	76.6	78.9	72.2	72.3	77.7	73.4	66.8	67.3
DeepSeek-R1-671B	84.8	86.3	81.3	81.7	87.2	86.8	81.8	82.6
Retrieval-augmented Reasoning								
Precedent-based-RAG-Qwen2.5-32B	88.5	88.2	85.8	85.7	89.4	87.2	83.7	84.5
Precedent-based-RAG-QwQ-32B	89.4	89.9	87.3	87.1	90.4	88.4	85.2	85.4
Precedent-based-RAG-Qwen2.5-72B	88.1	87.5	85.1	84.9	89.4	86.8	83.9	84.0
Search-o1-QwQ-32B	81.8	85.3	78.8	79.3	83.9	83.3	76.4	77.4
Agentic Retrieval-augmented Reasoning								
GLARE-Qwen2.5-32B(ours)	89.8	89.8	87.8	87.8	90.4	89.2	87.3	87.5
GLARE-QwQ-32B(ours)	89.7	90.7	88.6	88.6	91.3	90.6	88.3	88.5

Table 2: Performance comparison on CAIL 2018 dataset. The best results are in bold.

on large-scale Chinese legal corpora, which enhances its ability to process longer legal documents and capture complex contextual semantics.

LLM-based Methods: These methods utilize LLMs to perform legal reasoning in zero-shot or few-shot settings (Brown et al., 2020). **Direct Reasoning** directly feeds the case facts into the LLM to predict the applicable law articles and charges, without relying on any retrieval augmentation or additional external context. The models used in this setting include Qwen2.5-32B/72B-Instruct (Yang et al., 2024a), QwQ-32B (Team, 2025), and DeepSeek-R1-671B (Guo et al., 2025).

Retrieval-augmented Reasoning: (1) *Precedent-based RAG* enhances reasoning by retrieving top-5 precedents including their facts and labels, which are appended to the prompt. The models used in this setting include Qwen2.5-32B/72B-Instruct (Yang et al., 2024a), QwQ-32B (Team,

2025). (2) *Search-o1* (Liet al., 2025) dynamically retrieves external knowledge when it encounters uncertain or ambiguous knowledge in the general domain. We use reasoning model QwQ-32B (Team, 2025) in this setting.

basemodel to run the full reasoning pipeline. For generation, we set the following parameters: a maximum of 32,768 tokens and temperature of 0.6. For charge expansion, we set the top- k expanded charges to 3 in each expansion method. For precedent retrieval, we use SAILER (Li et al., 2023) to encode case facts and set the top- k retrieved precedents to 5. In the legal search module, we utilize SerperAPI² with the region configured for China and the number of returned results limited to the top 10. For charges that are not in the predefined label set, we map them to the most similar charge within the label set using BGE (Xiao et al., 2024).

Experiment Results

The results are reported in Table 2 and Appendix D, and next we will analyze the experimental results:

1. Our method has demonstrated consistent performance improvements in both charge pre-

diction and law article prediction tasks, highlighting the effectiveness of our agentic reasoning approach for LJP. Compared to the Direct Reasoning setting, our method improves charge prediction by 7.7% and law article prediction by 11.5% in F1 score. When compared with Retrieval-augmented Reasoning, it achieves an improvement of 1.5% on charge prediction and 3.1% on law article prediction in F1 score. In addition, our method not only performs well on the large reasoning models, but also effectively promotes the reasoning ability of the instruct models, indicating that our three mod-

ting, although the latter have significantly larger parameters sizes. The key reasons are as follows:

- (1) BERT frames charge and article prediction as multi-class classification tasks, enabling direct mapping from facts to fixed labels, which aligns well with the task. In contrast, LLMs take a generative approach and without legal-specific training they often fail to make accurate predictions.
- (2) Our dataset includes many rare and confusing charges. Fine-tuned BERT models trained on legal corpora can better distinguish these nuanced charges, while even large LLMs

²<https://serper.dev>

Methods	CAIL2018				CMDL			
	Charge		Law Article		Charge		Law Article	
	Acc.	Ma-F	Acc.	Ma-F	Acc.	Ma-F	Acc.	Ma-F
Direct Reasoning Qwen2.5-32B	60.2	39.3	63.7	41.2	57.4	64.7	57.9	63.7
QwQ-32B	78.4	57.0	79.2	58.2	67.9	72.9	69.8	74.2
Retrieval-augmented Reasoning								
Precedent-based-RAG-Qwen2.5-32B	82.6	62.7	83.0	62.3	65.7	69.5	65.3	67.6
Precedent-based-RAG-QwQ-32B	84.6	65.5	84.6	67.6	72.8	74.8	71.3	73.2
Agentic Retrieval-augmented Reasoning								
GLARE-Qwen2.5-32B(ours)	86.9	68.6	86.5	68.3	73.5	75.5	71.9	73.4
GLARE-QwQ-32B(ours)	90.7	75.7	91.1	75.4	76.0	79.5	74.0	76.7

Table3: Performance comparison on difficult charges.

ules effectively supplement legal knowledge and thereby enhance reasoning performance.

2. In contrast to direct reasoning, precedent-based RAG enhances prediction performance through precedent retrieval. Large reasoning models like QwQ-32B and DeepSeek-R1-671B outperform other instruct models in direct reasoning, indicating that LJP inherently requires multi-step reasoning and slow thinking. Precedent-based RAG improves performance across models of various scales by incorporating precedent retrieval. For example, QwQ-32B sees an 8.39% F1 improvement in law article prediction. However, precedent-based RAG only provides the case facts and labels of precedents, leading models to rely on similarity matching and copy judgment predictions rather than truly reason. Additionally, Search-o1 retrieves case facts which may introduce noises, rather than specific legal knowledge, thus underperforming compared to direct reasoning.

3. LLM-based methods outperform classification methods. However, BERT demonstrates superior performance compared to LLMs such as Qwen2.5-72B-Instruct in the direct reasoning set-

lack the domain knowledge needed for such difficult charges.

Ablation Study

Methods	Charge		Law Article	
	Acc.	Ma-P	Acc.	Ma-F
w/o CEM	89.6	87.7	90.3	85.2
w/o PRD	80.0	78.1	81.6	75.4
w/o LSAR	89.6	87.9	90.4	86.5
GLARE(ours)	89.7	88.6	91.3	88.5

Table4: Ablation Study. The best results are in bold.

To evaluate the effectiveness of each component in the GLARE framework, we conducted ablation experiments with the following strategies: (1) **w/o CEM**: The Charge Expansion Module is removed, so the model cannot expand a diverse set of candidate charges. (2) **w/o PRD**: The Precedents Reasoning Demonstration module is removed, so the model cannot refer to reasoning path from precedents. (3) **w/o LSAR**: The Legal Search-Augmented Reasoning module is removed, disabling the model's ability to supplement its knowl-

edge via external legal search when faced with ambiguous or unfamiliar charges.

As shown in Table 4, the removal of any single module results in degraded performance. In particular, removing PRD causes the most significant degradation: the accuracy of charge prediction drops from 89.7% to 80%. This highlights the crucial role of precedent reasoning path in enhancing legal judgment prediction. Removing CEM weakens the model's ability to recognize ambiguous or low-frequency charges, while LSAR helps the model fill knowledge gaps by retrieving authoritative legal information. Overall, the GLARE framework performs best across all metrics, validating the strength of agentic reasoning in legal judgment prediction.

Fact: The defendant, Song, purchased three general machine-printed invoices with a total face value of RMB 600,000, and subsequently submitted them for accounting purposes.

(1) **The overall inference efficiency is relatively high.** The average reasoning rounds for each case is 5.17 and the average call numbers for each module is between 1.7 and 1.8 times, indicating that the module scheduling is well-balanced, without obvious redundancy or repeated invocation. So the overall delay is within our acceptable range.

(2) **The CEM module is the most efficient.** In both expansion methods, the charge structures are established offline in advance, so its computation cost is low and runtime is minimal. As shown in the figure 5(b), compared with direct reasoning and RAG approaches, our method considers a comprehensive set of charges and performs a more thorough analysis.

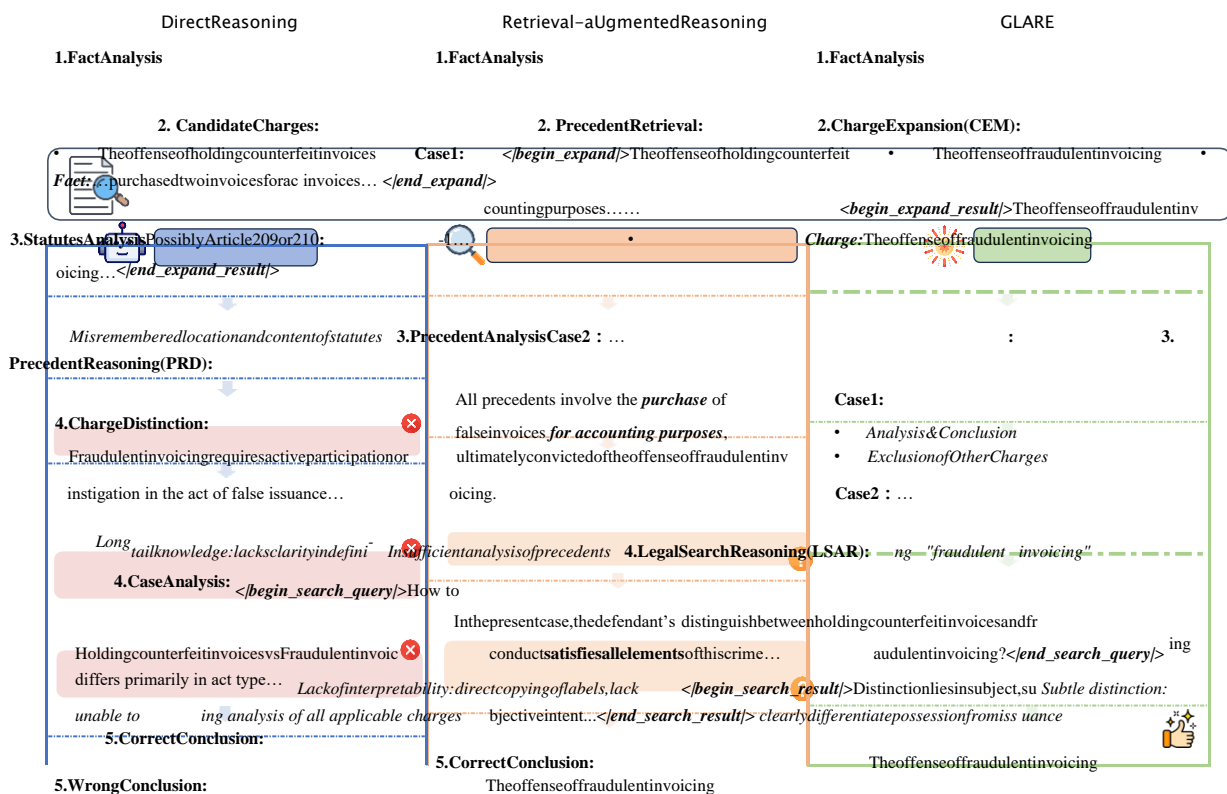


Figure 4: Case Study. The red part highlights the model's limitations due to insufficient internal knowledge, while the yellow part demonstrates the lack of interpretability in vanilla precedent-based RAG reasoning.

Efficiency Analysis

In this work, we focus on multi-step reasoning and slow thinking for legal judgment prediction, so the latency is less important. Nevertheless, we still conducted an analysis to further understand each module. Based on the analysis of Figure 5(a), we can draw the following conclusions:

(3) **The PRD module has the highest latency but within an acceptable range.** Since this module needs to encode the entire case description and the case text is usually long, the reasoning time is relatively long. However, the PPR module can provide the reasoning path of precedents and has significant reasoning interpretability.

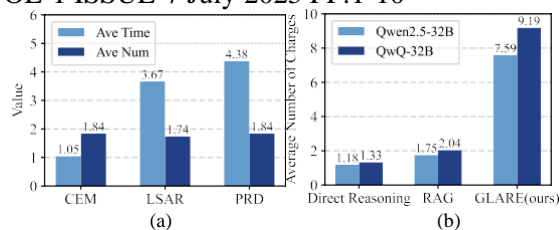


Figure 5: (a) Efficiency analysis of each module.(b) Average charge numbers of different methods.

Case Study

As shown in the Figure 4, we conducted case study on three LLM-based methods to further verify the effectiveness and interpretability of our method. Direct Reasoning relies on the LLM's internal knowledge, which may be inaccurate or insufficient, leading to incorrect judgments. RAG methods often lack explicit links between retrieved cases and final decisions, making it hard to trace how external knowledge affects reasoning. However, our method ensures that each reasoning step has a clear knowledge basis through the explicit invocation of three modules, thereby extending the reasoning chain.

Performance on Difficult

Charges

To evaluate GLARE's ability to handle challenging charges requiring long-tail knowledge, We conducted experiments on low-frequency charges with less than 100 cases (e.g., the crime of bribing a unit) and confusing charges (e.g., robbery vs. snatching). The results are reported in Table 3, which reveal two key insights: (1) Our method dynamically acquires critical legal knowledge, outperforming Direct Reasoning by over 10% and Retrieval-augmented Reasoning by over 5%. (2) RAG-based methods struggle to retrieve relevant precedents for such charges, leading to poor performance, while direct reasoning fall short due to limited long-tail knowledge. These results highlight the strength of our external modules in supplementing legal reasoning with critical knowledge.

5 Conclusion

In this study, we propose a novel framework, GLARE, to address the legal gaps in legal reasoning. GLARE dynamically acquires key legal knowledge to improve the breadth and depth of reasoning. Experimental results demonstrate the effectiveness of our approach, which not only improves prediction performance but also generates complete reasoning chains that enhance the inter-

pretability of LJP tasks. We believe that GLARE holds great potential for real-world legal applications and will contribute meaningfully to the advancement of intelligent judicial systems.

Limitations

Generalizability We adopted the legal dataset from China Judgments Online to verify the applicability of the method in the China's judicial system. However, the GLARE framework is applicable to countries following both common law and civil law systems. When applied to the actual judicial practice of a specific country, we need to inject the specific legal knowledge base of each country and adapt to the local judicial culture.

Efficiency Our method promotes the reasoning ability of the model through multiple rounds of reasoning and the invocation of three modules. Although this process has an increased time cost compared to the traditional direct reasoning method, the task of legal judgment prediction itself is a task that requires multi-step reasoning and slow thinking. Moreover, this time cost is much less than the time needed for humans to analyze cases in real life. Therefore, such a time cost is acceptable.

Ethical Discussion

Potential Bias in Legal Data Large language models may learn historical bias from legal judgments in training data. In practice, judicial decisions may be influenced by many external factors, such as public opinion, regional differences or the personal inclination of judges. We need to identify possible biases before deploying such models in real-world scenarios.

Human-Centric Deployment Our system is designed to assist judges by providing supplementary recommendations rather than replacing human decision-making. We advise users to critically evaluate the model's predictions and make independent decisions about their adoption, rather than uncritically accepting the model's reasoning.

Artificial Intelligence and Law , 28(2), 237-266.

References:

1. Ashley, K. D. (2017). *Artificial Intelligence and Legal Analytics: New Tools for Law Practice in the Digital Age* . Cambridge University Press.
2. Susskind, R. (2019). *Online Courts and the Future of Justice* . Oxford University Press.
3. Zeng, J., & Wang, T. (2023). A Survey of Legal Judgment Prediction: Datasets, Metrics, Models, and Challenges. *AI Open* , 4, 1-12. Why: A recent survey that provides a comprehensive overview of the LJP field, its state-of-the-art, and open problems.
4. Medvedeva, M., Vols, M., & Wieling, M. (2020). Using machine learning to predict decisions of the European Court of Human Rights.
5. Zhong, H., Guo, Z., Tu, C., Xiao, C., Liu, Z., & Sun, M. (2018). Legal judgment prediction via topological learning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing* (pp. 3540-3549).
6. Chalkidis, I., Kampas, D., & Androutsopoulos, I. (2019). Neural legal judgment prediction in English. In
7. Xu, N., Wang, P., Chen, L., Pan, L., Wang, X., & Zhao, J. (2020). Distinguish Confusing Law Articles for Legal Judgment Prediction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics* (pp. 3086-3095).
8. Katz, D. M., Bommarito II, M. J., & Blackman, J. (2017). A general approach for predicting the behavior of the supreme court of the united states. *PloS one* , 12(4), e0174698.
9. Savelka, J., Ashley, K. D., Gray, M. A., & Walker, V. R. (2023). Explaining Legal Concepts with Augmented Large Language Models (GPT-4). *arXiv preprint arXiv:2306.09525* .
10. Wooldridge, M. (2009). *An Introduction to MultiAgent Systems* . John Wiley & Sons.
11. Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., ... & Wen, J. (2023). A Survey on Large Language Model based Autonomous Agents. *arXiv preprint arXiv:2308.11432* . Why: An excellent and recent survey on the very hot topic of using LLMs as the "brains" of autonomous agents, directly relevant to your approach.
12. Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, K., ... & Tang, J. (2023). The Rise and Potential of Large Language Model Based Agents: A Survey. *arXiv preprint arXiv:2309.07864* .
13. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., ... & Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large

14. Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., & Cao, Y. (2022). ReAct: Synergizing Reasoning and Acting in Language Models. *arXiv preprint arXiv:2210.03629* .

15. Long, J. (2023). Large Language Model Guided Tree-of-Thought. *arXiv preprint arXiv:2305.08291* .

16. Surden, H. (2019). Artificial Intelligence and Law: An Overview. *Georgia State University Law Review* , 35(4).

17. Citron, D. K., & Pasquale, F. (2014). The Scored Society: Due Process for Automated Predictions. *Washington Law Review* , 89, 1. Why: A classic legal article discussing the due process requirements for automated decisionmaking systems, crucial for any work on "judgment prediction."

18. Goodman, B., & Flaxman, S. (2017). European Union regulations on algorithmic decision-making and a "right to explanation". *AI Magazine* , 38(3), 50-57.

19. Felzmann, H., Fosch-Villaronga, E., Lutz, C., & Tamò-Larrieux, A. (2020). Transparency you can trust: Transparency requirements for artificial intelligence between legal norms and contextual concerns. *Big Data & Society* , 7(1).

20. Yeung, K. (2017). 'Hypernudge': Big Data as a mode of regulation by design. *Information, Communication & Society* , 20(1), 118-136.