# DianJin-R1: Evaluating and Enhancing Financial Reasoning in Large Language Models

**Jie Zhu**[1], **Qian Chen**[1], **Huaixia Dou**[1,2], **Junhui Li**[2],
**Lifan Guo**[1], **Feng Chen**[1], **Chi Zhang**[1]
[1]Qwen DianJin Team, Alibaba Cloud Computing
[2]Soochow University

🙂 https://huggingface.co/DianJin
[icon] https://modelscope.cn/organization/tongyi_dianjin
⚙ https://github.com/aliyun/qwen-dianjin
[icon] https://tongyi.aliyun.com/dianjin

## Abstract

Effective reasoning remains a core challenge for large language models (LLMs) in the financial domain, where tasks often require domain-specific knowledge, precise numerical calculations, and strict adherence to compliance rules. We propose DianJin-R1, a reasoning-enhanced framework designed to address these challenges through reasoning-augmented supervision and reinforcement learning. Central to our approach is DianJin-R1-Data, a high-quality dataset constructed from CFLUE, FinQA, and a proprietary compliance corpus (Chinese Compliance Check, CCC), combining diverse financial reasoning scenarios with verified annotations. Our models, DianJin-R1-7B and DianJin-R1-32B, are fine-tuned from Qwen2.5-7B-Instruct and Qwen2.5-32B-Instruct using a structured format that generates both reasoning steps and final answers. To further refine reasoning quality, we apply Group Relative Policy Optimization (GRPO), a reinforcement learning method that incorporates dual reward signals: one encouraging structured outputs and another rewarding answer correctness. We evaluate our models on five benchmarks: three financial datasets (CFLUE, FinQA, and CCC) and two general reasoning benchmarks (MATH-500 and GPQA-Diamond). Experimental results show that DianJin-R1 models consistently outperform their non-reasoning counterparts, especially on complex financial tasks. Moreover, on the real-world CCC dataset, our single-call reasoning models match or even surpass the performance of multi-agent systems that require significantly more computational cost. These findings demonstrate the effectiveness of DianJin-R1 in enhancing financial reasoning through structured supervision and reward-aligned learning, offering a scalable and practical solution for real-world applications.

## 1 Introduction

Recent advances in large language models (LLMs) have led to growing interest in enhancing their reasoning abilities. Models such as OpenAI o1 (OpenAI, 2024), DeepSeek R1 (Guo et al., 2025) and QwQ (Qwen, 2024) have shown that explicitly modeling reasoning processes can significantly boost performance on complex tasks (Zhong et al., 2024). Despite these improvements, recent evaluations on financial benchmarks (Xie et al., 2023; 2024; Zhu et al., 2024; Chen et al., 2024; Qian et al., 2025; Liu et al., 2025) reveal that reasoning in this domain remains particularly challenging, given the need for domain-specific knowledge, accurate numerical reasoning, and strict compliance with regulatory requirements. Effectively addressing these challenges calls for specialized reasoning strategies capable of handling both structured financial information and open-ended problem solving. In re-

arXiv:2504.15716v1 [cs.AI] 22 Apr 2025

sponse, we introduce DianJin-R1, LLMs that incorporate reasoning-augmented supervision and reinforcement learning to enhance performance on financial reasoning tasks.

We begin by constructing a high-quality reasoning dataset, DianJin-R1-Data, using three major sources: CFLUE (Zhu et al., 2024), FinQA (Chen et al., 2021), and our proprietary compliance dataset for the task of Chinese Compliance Check (CCC). CFLUE, which includes over 31,000 reasoning-annotated multiple-choice and open-ended questions from financial qualification mock exams, plays a central role in training due to its scale and diversity. FinQA provides numerical reasoning questions, while CCC focuses on complex compliance scenarios requiring multi-step logic. To ensure the quality of reasoning, we adopt a verification process using GPT-4o (OpenAI, 2024) to check for alignment between generated answers, reasoning steps, and reference explanations. This process results in a reliable set of reasoning-augmented and non-reasoning samples, supporting more robust model training.

For supervised fine-tuning (SFT), we train DianJin-R1-7B and DianJin-R1-32B, based on Qwen2.5-7B-Instruct and Qwen2.5-32B-Instruct (Yang et al., 2024), to generate both the reasoning process and final answers using a structured output format with `<think>` and `<answer>` tags. To further improve reasoning quality, we apply Group Relative Policy Optimization (GRPO) (Shao et al., 2024), a reinforcement learning algorithm that introduces two reward signals: a format reward to encourage structured outputs and an accuracy reward to promote answer correctness. These mechanisms guide the model to produce coherent, verifiable reasoning paths and reliable answers.

We evaluate our DianJin-R1 models, along with other general reasoning and non-reasoning models, across a diverse set of benchmarks, including CFLUE, FinQA, CCC, MATH-500 (Hendrycks et al., 2021), and GPQA-Diamond (Rein et al., 2024). The results demonstrate that reasoning-augmented models consistently outperform their non-reasoning counterparts, especially in the financial domain. Notably, training on CFLUE alone yields substantial gains across all tasks, and combining all datasets further enhances performance. Our analysis also highlights the benefit of reinforcement learning, particularly when the reward signals align with the task domain.

Finally, we demonstrate a practical application of our approach on the CCC dataset, where a multi-agent system based on LLMs is employed to perform condition-based compliance checks. By assigning specialized agents to each decision node in the workflow, the system effectively integrates intermediate reasoning steps to arrive at the final compliance judgment.

In summary, DianJin-R1 presents a scalable and effective strategy for enhancing financial reasoning in LLMs by combining high-quality supervision, structured reasoning generation, and reward-driven refinement through reinforcement learning.

## 2 DianJin-R1-Data Construction

### 2.1 Data Source

Our dataset originates from different sources: two open-source datasets and an in-house dataset.

**CFLUE (Zhu et al., 2024).** It is an open-source Chinese benchmark designed to assess the performance of LLMs on a variety of natural language processing (NLP) tasks within the financial domain. Its knowledge assessment component includes 38,638 multiple-choice financial exam questions, sourced from 15 types of financial qualification mock exams that cover various subjects and difficulty levels. To construct a high-quality subset for our study, we apply a three-step filtering process—focusing on question length, difficulty, and ambiguity. First, we apply a length filter to remove questions with fewer than 15 tokens, as these typically require minimal reasoning and offer limited value for assessing deeper understanding. Second, since simple QA pairs may not significantly enhance reasoning ability (Ye et al., 2025; Muennighoff et al., 2025), we apply a difficulty filter to discard questions that are correctly answered by all smaller language models, including

LLaMA-3.1-8B and Qwen2.5-7B-Instruct. This helps ensure that the remaining questions demand deeper reasoning. Finally, in the ambiguity filter, we use GPT-4o (OpenAI, 2023) to eliminate questions that contain ambiguous wording, thereby ensuring each question is clear and well-defined. Through this filtering process, we curate a refined set of high-quality multiple-choice questions, each with an unambiguous correct answer, suitable for evaluating financial reasoning in LLMs. Notably, most of these question-answer pairs are accompanied by detailed explanations, which serve as valuable reasoning annotations.

**FinQA (Chen et al., 2021).** FinQA is an open-source English benchmark containing 8,281 financial question-answer pairs that require numerical reasoning over financial reports. For our study, we apply the same length and difficulty filters as used in the CFLUE dataset to ensure quality and complexity. As a result, we curate a high-quality subset of QA pairs, well-suited for evaluating financial reasoning in English-language contexts.

**CCC (Chinese Compliance Check).** It is an in-house dataset designed to detect compliance violations by service agents in real-world Chinese financial customer-service dialogues. Each instance in CCC includes the full conversation between a customer and a service agent, along with associated meta-information such as customer service actions (e.g., ticket escalation), and other relevant contextual data. The data is sourced from an online quality inspection system used in actual customer service operations, and each instance has been manually reviewed to ensure labeling accuracy. While compliance checking can be framed as a binary classification task (violation vs. non-violation), it is inherently more complex, as it requires evaluating whether the service agent's behavior adheres to a set of domain-specific regulatory guidelines. Specifically, we sample from the manually validated data, ensuring a roughly balanced distribution between compliant and non-compliant cases.

## 2.2 Reasoning Dataset Construction

Considering the differences between the dataset CCC and the two others CFLUE and FinQA, we adopt different methods for reasoning construction for multiple-choice questions in CFLUE (Section 2.2.1), QA pairs of FinQA (Section 2.2.2) and dialogues of CCC (Section 2.2.3). Table 1 shows the statistics.

### 2.2.1 Reasoning Generation for CFLUE Questions

We denote our selected multiple-choice questions from CFLUE as $D_{\text{CFLUE}_{\text{MCQ}}} = (x_i, e_i, y_i)$, where each $x_i$ includes both the question and its corresponding set of multiple-choice options, $e_i$ is the detailed explanation, and $y_i$ includes both a summarized reasoning process and the correct answer, typically indicated by a letter such as $A$, $B$, $C$, or $D$.

**Open-ended question generation.** We begin by using GPT-4o to convert each multiple-choice question in $D_{\text{CFLUE}_{\text{MCQ}}}$ into an open-ended format, where $x$ denotes the input question and $y$ represents the correct answer. Figure 3 in Appendix A provides an illustration of this conversion process. We denote the resulting dataset of open-ended questions as $D_{\text{CFLUE}_{\text{OE}}} = \{(x_i, y_i)\}$.

**Reasoning generation.** For each pair $(x_i, y_i)$ in $D_{\text{CFLUE}_{\text{MCQ}}}$, we leverage DeepSeek-R1 (Guo et al., 2025), a model known for its strong reasoning capabilities, to generate a chain-of-thought (CoT) $r_i$ along with a predicted answer $a_i$ as follows:

$$r_i, a_i = \text{DS}(x_i). \tag{1}$$

Next, we employ GPT-4o as a verifier to assess two key aspects of the generated output: (1) whether the predicted answer $a_i$ matches the gold answer $y_i$, and (2) whether the generated reasoning $r_i$ is consistent with the reference explanation $e_i$. If both conditions are satisfied, we retain the instance $(x_i, r_i, y_i)$ as a valid reasoning sample. If not, we retry the reasoning generation process, allowing up to $T$ attempts (with $T = 3$ in this paper). If all attempts fail

to produce a correct answer and consistent reasoning, the instance $(x_i, y_i)$ is considered a hard case and preserved as a non-reasoning sample. We denote the resulting reasoning-augmented dataset as $R_{\text{CFLUE}_{\text{MQC}}} = \{(x_i, r_i, y_i)\}$ and the hard, non-reasoning dataset as $G_{\text{CFLUE}_{\text{MQC}}} = \{(x_i, y_i)\}$.

We apply the same procedure to the open-ended question dataset $D_{\text{CFLUE}_{\text{OE}}}$, producing the corresponding reasoning-augmented dataset $R_{\text{CFLUE}_{\text{OE}}} = \{(x_i, r_i, y_i)\}$ and the hard, non-reasoning dataset as $G_{\text{CFLUE}_{\text{OE}}} = \{(x_i, y_i)\}$.

### 2.2.2 Reasoning Generation for FinQA Questions

Unlike the instances in CFLUE, the QA pairs in FinQA are already in an open-ended format. We denote the FinQA dataset similarly as $D_{\text{FinQA}} = \{(x_i, y_i)\}$.

We then apply the procedure of reasoning generation as open-ended questions in CFLUE to the QA pairs in $D_{\text{FinQA}}$. As a result, we obtain the reasoning-augmented dataset from FinQA as $R_{\text{FinQA}} = \{(x_i, r_i, y_i)\}$ and the hard, non-reasoning dataset as $G_{\text{FinQA}} = \{(x_i, y_i)\}$.

### 2.2.3 Reasoning Generation via Multi-Agent for CCC Dialogues

We denote the CCC dataset as $D_{\text{CCC}} = \{(x_i, y_i)\}$, where $x_i$ is a dialogue and $y_i$ is the corresponding answer, which provides a summarized reasoning process and a final conclusion on whether a compliance violation has occurred. It is challenging for LLMs to directly generate reasoning from a dialogue $x_i$, since even humans typically rely on a set of guidelines to determine whether a service agent has violated compliance. To replicate this process, we have developed a workflow based on these guidelines, outlining the detailed steps for identifying compliance violations. The workflow begins with a start node and ends with two possible outcome nodes. One outcome indicates that the dialogue has no compliance violation, while the other indicates a compliance violation. All other nodes in the workflow are internal condition nodes, each evaluating whether a specific condition is met and triggering corresponding actions based on the result.

Although using a multi-agent LLM-based system to identify compliance violations is a natural approach, it significantly increases inference costs due to multiple agent interactions. As an alternative, we use a multi-agent LLM system to generate reasoning. Specifically, we employ an LLM-based agent (i.e., `Qwen2.5-72B-Instruct` in this study) for each condition node. For each pair $(x_i, y_i)$ in $D_{\text{CCC}}$, we strictly follow the workflow and ask the corresponding agents to generate intermediate CoTs and the associated intermediate answers. The intermediate CoTs and intermediate answers are denoted as $(r_i^0, a_i^0, \cdots, r_i^{n_i}, a_i^{n_i})$, where $n_i$ is the number of condition nodes involved before reaching the outcome node. The final answer, $a_i$, is determined by the outcome node.

If the generated answer $a_i$ matches the gold answer $y_i$, we use `GPT-4o` to merge all the intermediate CoTs $(r_i^0, \cdots, r_i^{n_i})$ into a final, unified CoT $r_i$ and retain the instance $(x_i, r_i, y_i)$. If the answer does not match, we retry the reasoning generation process up to $T$ attempts. Similarly, we denote the resulting reasoning-augmented dataset as $R_{\text{CC}} = \{(x_i, r_i, y_i)\}$ and the hard, non-reasoning dataset as $G_{\text{CC}} = \{(x_i, y_i)\}$. Figure 1 presents an example of how the final unified CoT is generated from the intermediate CoTs produced by the multi-agent system.

## 2.3 Model Training

In this section, we present the details for training LLMs to perform financial reasoning. As shown in Figure 2, the training consists of two stages: learning reasoning with supervised fine-tuning (SFT) and enhancing reasoning with reinforcement learning (RL).

### 2.3.1 Learning Reasoning with SFT

The reasoning datasets—$R_{\text{CFLUE}_{\text{MCQ}}}$, $R_{\text{CFLUE}_{\text{OE}}}$, $R_{\text{FinQA}}$, and $R_{\text{CCC}}$—are utilized to fine-tune LLMs to generate a chain-of-thought (CoT) followed by a final answer. Each training in-

Figure 1: An example of reasoning data synthesized by a multi-agent system.

| Dataset | Language | Size | $Q_{token}$ | $R_{token}$ | $A_{token}$ |
|---------|----------|------|-------------|-------------|-------------|
| Used in SFT | | | | | |
| CFLUE$_{MCQ}$ | Chinese | 26,672 | 134.85 | 807.42 | 95.71 |
| CFLUE$_{OE}$ | Chinese | 5,045 | 49.28 | 857.04 | 485.60 |
| FinQA | English | 4,851 | 1048.38 | 1576.91 | 148.42 |
| CCC | Chinese | 1,800 | 1695.78 | 884.29 | 69.64 |
| Used in RL | | | | | |
| CFLUE$_{MCQ}$ | Chinese | 4096 | 132.40 | - | 2.15 |

Table 1: Overview of datasets used in DianJin-R1-Data.

stance $(x, r, y)$ consists of a question $x$, a reasoning path $r$ (formatted as <think>· · ·</think>), and an answer $y$ (formatted as <answer>· · ·</answer>). During fine-tuning, the question $x$ serves as the input to the model, while the reasoning $r$ and final answer $y$ are treated as the target output, enabling the model to learn to produce coherent reasoning steps along with the correct solution.

### 2.3.2   Enhancing Reasoning with RL

The hard and non-reasoning dataset $G_{\text{CFLUE}_{\text{MCQ}}}$ is used in the reinforcement learning (RL), which aims to further enhance the reasoning skills.[1]  Specifically, we adopt the Group Relative Policy Optimization (GRPO) algorithm (Shao et al., 2024) for RL, incorporating two reward mechanisms: a format reward to ensure the generated output adheres to the desired structure, and an accuracy reward to encourage correct answers.

---

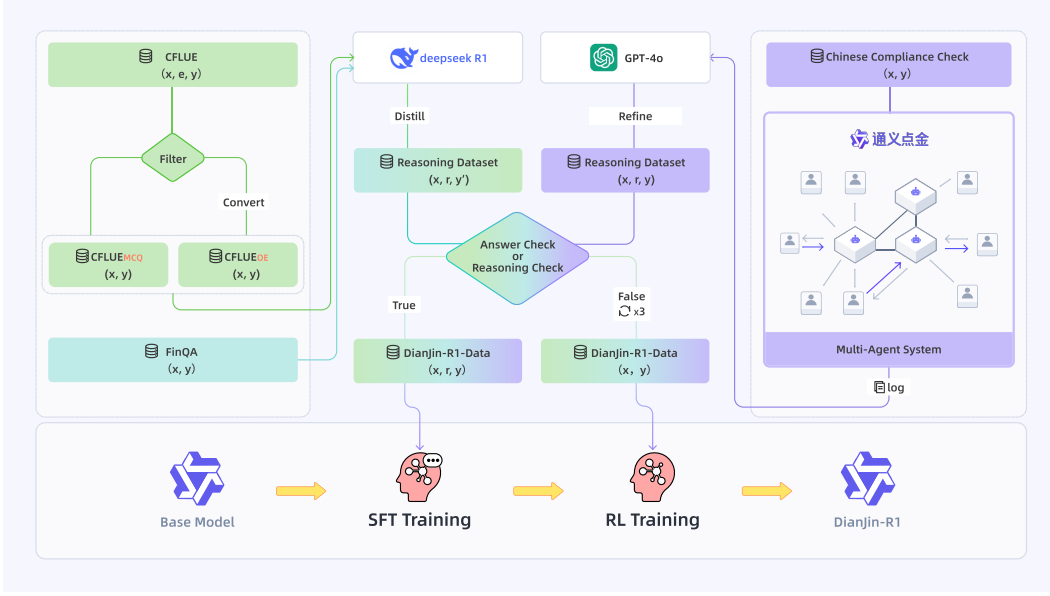[1]In the future, we will include $G_{\text{CFLUE}_{\text{OE}}}$, $G_{\text{FinQA}}$, and $G_{\text{CCC}}$ in RL.

Figure 2: Illustration of two-step training for DianJin-R1.

- Format reward: We incorporate a format reward to promote well-structured outputs. For a given response $y$, a reward score of 1 is granted if the output contains exactly one reasoning segment enclosed within `<think>`$\cdots$`</think>` tags and one final answer enclosed within `<answer>`$\cdots$`</answer>` tags, with no additional content outside these boundaries. If this strict formatting criterion is not met, the output receives a reward score of 0.

- Accuracy reward: We use an accuracy reward to promote answer correctness. Specifically, for a given output $y$, if the content enclosed within the `<answer>`$\cdots$`</answer>` tags exactly matches the reference answer,[2] the model receives a reward score of 1; otherwise, it receives a score of 0. This mechanism encourages the generation of accurate and reliable final answers.

## 3 Experimentation

### 3.1 Experimental Setups

**Training Data.** Table 1 summarizes the statistics for DianJin-R1-Data. Since CFLUE$_{\text{MCQ}}$ constitutes a large portion of the SFT data, we shuffle it together with the other datasets to prevent overfitting to any single source.

**Model Training.** Using the constructed DianJin-R1-Data and our proposed training method, we train DianJin-R1-7B and DianJin-R1-32B based on Qwen2.5-7B-Instruct and Qwen2.5-32B-Instruct, respectively. Our model training consists of two stages,conducted on NVIDIA A100 GPUs. In the SFT stage, the 7B model is trained on a single node with 8 GPUs, and the 32B model uses 32 GPUs across 4 nodes. We leverage DeepSpeed's Zero-3 for optimization, using a learning rate of $1.0 \times 10^{-5}$, a sequence length of 16K, and bf16 precision over 3 epochs. We apply gradient accumulation with 16 steps to simulate a larger batch size. In the RL stage, we perform 8 rollouts per sample, with a train batch size of 1024 and a rollout batch size of 256. We use a learning rate of $1.0 \times 10^{-6}$ and a sampling temperature of 0.6 over 5 epochs to balance learning dynamics.

---

[2]In the current version, the accuracy reward is applied exclusively to multiple-choice questions from $R_{\text{CFLUE}_{\text{MCQ}}}$.

**Evaluation Datasets.** We evaluate our models on three financial benchmarks, including the test sets of CFLUE, FinQA, and our proprietary in-house dataset, CCC. To further assess general reasoning capabilities, we also include two widely-used benchmarks: MATH-500 (Hendrycks et al., 2021) and GPQA-Diamond (Rein et al., 2024). For each dataset, we report the accuracy—defined as the proportion of correctly answered questions—and compute the average accuracy across all test sets. Among them, CFLUE and CCC are Chinese-language datasets, while the others are in English. The detailed statistics of these test sets are summarized in Table 2. For FinQA and CCC, we use GPT-4o to evaluate the correctness of each answer, following the prompts shown in Figure 8 and Figure 9 in Appendix B. For the other test sets, we extract the predicted answers using rule-based methods and compare them directly with the gold answers.

| Dataset | Language | Size |
|---|---|---|
| Financial domain | | |
| CFLUE | Chinese | 3,864 |
| FinQA | English | 1,147 |
| CCC | Chinese | 200 |
| General domain | | |
| MATH-500 | English | 500 |
| GPQA-Diamond | English | 198 |

Table 2: Statistics of the test sets.

**Baselines** We compare our models against two categories of LLMs. The first includes general LLMs without explicit reasoning capabilities: GPT-4o (OpenAI, 2024), DeepSeek-V3 (Liu et al., 2024), Qwen2.5-7B-Instruct (Yang et al., 2024), Qwen2.5-32B-Instruct, and Qwen2.5-72B-Instruct. The second category consists of general LLMs equipped with reasoning abilities, including DeepSeek-R1 (Guo et al., 2025), DeepSeek-R1-Distill-Qwen-7B, DeepSeek-R1-Distill-Qwen-14B, DeepSeek-R1-Distill-Qwen-32B, and QwQ-32B (Qwen, 2024).

## 3.2 Experimental Results

Table 3 compares the performance of LLMs with and without explicit reasoning. Overall, the results demonstrate that models incorporating reasoning generally outperform their non-reasoning counterparts, with the exception of DeepSeek-R1-Distill-Qwen-7B. The results reveal the following key points:

- On the three financial test sets (i.e., CFLUE, FinQA, and CCC), our DianJin-R1 models significantly outperform the base models (Qwen2.5-7B-Instruct and Qwen2.5-32B-Instruct), especially on CFLUE and CCC. For instance, DianJin-R1-32B improves accuracy from 77.95 to 86.74 on CFLUE, from 79.51 to 80.82 on FinQA, and from 56.50 to 96.00 on CCC. Encouragingly, DianJin-R1-32B achieves the highest accuracy on these financial tasks, surpassing even the performance of DeepSeek-R1, a model known for its strong reasoning capabilities.

- On the two general-domain test sets (i.e., MATH-500 and GPQA-Diamond), we also observe performance improvements in the DianJin-R1 models compared to their respective base models. For example, DianJin-R1-32B improves accuracy from 81.00 to 88.20 on MATH-500, and from 44.95 to 58.59 on GPQA-Diamond. This suggests that training on financial reasoning data can enhance general reasoning capabilities to some extent. However, since our SFT and RL training pipelines do not incorporate any general-domain reasoning datasets, the performance of DianJin-R1 models on these benchmarks remains lower than that of models with larger parameter sizes or those fine-tuned on general reasoning data.

- Although general-purpose reasoning models such as DeepSeek-R1 and QwQ-32B significantly improve performance on general reasoning benchmarks like MATH-500 and GPQA-Diamond, they do not always yield better results on financial

benchmarks. For instance, while DeepSeek R1 performs better than DeepSeek-V3 on CFLUE and CCC, it actually leads to a drop in performance on FinQA. Similarly, DeepSeek-R1-Distill-Qwen-7B performs worse than Qwen-2.5-7B-Instruct across all financial test sets. These findings are similar to those previously reported in Fino1 (Qian et al., 2025).

| Model | Financial | | | General | | Avg. |
|---|---|---|---|---|---|---|
| | CFLUE | FinQA | CCC | MATH | GPQA | |
| General models without explicit reasoning | | | | | | |
| GPT-4o | 71.68 | 79.16 | 50.00 | 77.93 | 39.56 | 63.67 |
| DeepSeek-V3 | 75.14 | **81.34** | 57.50 | 87.20 | 45.45 | 68.33 |
| Qwen2.5-7B-Instruct | 69.37 | 66.70 | 55.00 | 71.40 | 33.84 | 59.26 |
| Qwen2.5-32B-Instruct | 77.95 | 79.51 | 56.50 | 81.00 | 44.95 | 67.98 |
| Qwen2.5-72B-Instruct | 79.46 | 77.94 | 55.50 | 82.20 | 39.90 | 67.00 |
| General models with reasoning | | | | | | |
| DeepSeek-R1 | <u>86.64</u> | 79.81 | 67.50 | <u>94.80</u> | **66.16** | <u>78.98</u> |
| DeepSeek-R1-Distll-Qwen-7B | 48.39 | 66.09 | 41.50 | 90.20 | 45.96 | 58.43 |
| DeepSeek-R1-Distll-Qwen-14B | 70.83 | 76.63 | 50.00 | 93.20 | 54.55 | 69.04 |
| DeepSeek-R1-Distll-Qwen-32B | 78.52 | 77.00 | 52.00 | **95.00** | <u>63.64</u> | 73.23 |
| QwQ-32B | 83.49 | 78.38 | 52.00 | **95.00** | <u>63.64</u> | 74.50 |
| DianJin-R1 with reasoning | | | | | | |
| DianJin-R1-7B | 80.32 | 77.72 | <u>94.50</u> | 76.60 | 37.54 | 73.34 |
| DianJin-R1-32B | **86.74** | <u>80.82</u> | **96.00** | 88.20 | 58.59 | **82.07** |

Table 3: Performance comparison in accuracy across different test sets. Scores in **bold** and <u>underlined</u> indicate the best and second-best results, respectively.

## 3.3 Discussion

In this section, we use `Qwen2.5-7B-Instruct` as the backbone model to investigate the impact of RL and different datasets used during SFT.

**Effect of RL.** As shown in Table 4, SFT significantly enhances performance across all datasets. This indicates that SFT effectively equips the language models with reasoning capabilities. Furthermore, RL brings additional improvements on all datasets except FinQA. We suspect this exception may be due to the fact that all instances used for RL are in Chinese and sourced from CFLUE, whereas FinQA is in English. We plan to explore this further by incorporating English examples into the RL stage in future work.

| Learning | Financial | | | General | | Avg. |
|---|---|---|---|---|---|---|
| | CFLUE | FinQA | CCC | MATH | GPQA | |
| - | 69.37 | 66.70 | 55.00 | 71.40 | 33.84 | 59.26 |
| SFT | 75.48 | **80.28** | 93.00 | 74.87 | 34.85 | 71.70 |
| SFT + RL | **80.32** | 77.72 | **94.50** | **76.60** | **37.54** | **73.34** |

Table 4: Performance comparison in accuracy.

**Effect of different datasets used in SFT.** We use three data sources for SFT: CFLUE, FinQA, and CCC. To understand the individual contribution of each dataset, we evaluate model performance using different combinations of these sources. Table 5 presents the results. Among the datasets, CFLUE (which includes both CFLUE$_{MCQ}$ and CFLUE$_{OE}$) has the greatest impact. This is largely due to its large scale, as it contains more than 31,000 reasoning instances. When used alone, CFLUE significantly improves performance across all test sets, increasing overall accuracy from 59.26 to 65.67. Adding FinQA or CCC on top of CFLUE mainly enhances performance on their respective test sets, with limited influence

on the others. Finally, using all three datasets together during SFT results in the best overall performance.

| Dataset | Financial | | | General | | Avg. |
|---|---|---|---|---|---|---|
| | CFLUE | FinQA | CCC | MATH | GPQA | |
| - | 69.37 | 66.70 | 55.00 | 71.40 | 33.84 | 59.26 |
| CFLUE | 75.09 | 71.65 | 70.50 | 75.07 | **36.03** | 65.67 |
| CFLUE + FinQA | 75.20 | **80.85** | 70.00 | 75.33 | **36.03** | 67.48 |
| CFLUE + CCC | 75.16 | 71.78 | **95.00** | **76.27** | 35.35 | 70.71 |
| CFLUE + CCC + FinQA | **75.48** | 80.28 | 93.00 | 74.87 | 34.85 | **71.70** |

Table 5: Performance comparison in accuracy when different datasets are used in SFT. Here, RL is not applied.

**Effect of different systems on CCC.** As discussed earlier, it is natural to build a multi-agent LLM-based system to detect compliance violations in the CCC dataset. To achieve this, we follow the same reasoning generation workflow described in Section 2.2.3. Specifically, we assign an LLM agent (i.e., Qwen2.5-72B-Instruct), to each condition node in the workflow. Each agent is responsible for predicting the intermediate outcome of its corresponding condition. The final decision is then derived from the result at the outcome node based on the collective outputs of all agents. Table 6 presents the performance comparison on the CCC test set. Among the non-reasoning systems, introducing the multi-agent approach alone significantly improves accuracy from 55.50 to 95.00, highlighting the effectiveness of structured, condition-based reasoning for this task. However, this comes at a high cost, as it requires an average of 8.15 API calls per instance. Notably, our reasoning-based models, DianJin-R1-7B and DianJin-R1-32B, achieve comparable or even superior performance with just a single API call despite their smaller sizes. This demonstrates the strength of our approach in learning to reason effectively and organizing reasoning paths to handle complex compliance evaluations.

| System | Accuracy | #Calls |
|---|---|---|
| Qwen2.5-72B-Instruct | 55.50 | 1 |
| Qwen2.5-72B-Instruct (Multi-Agent) | 95.00 | 8.15 |
| DianJin-R1-7B | 94.50 | 1 |
| DianJin-R1-32B | **96.00** | 1 |

Table 6: Performance comparison in accuracy on the test set of CCC. #Calls denotes the average number of API calls required per test instance.

## 4   Conclusion and Future Work

We have presented DianJin-R1, a reasoning-augmented framework for large language models in the financial domain. This framework combines structured supervision with a reinforcement learning algorithm (GRPO) to enhance model performance on complex financial and compliance-related tasks. Through extensive experiments on diverse benchmarks and a real-world compliance system, we demonstrate that reasoning-aware training significantly improves both accuracy and interpretability.

In future work, we plan to explore alternative reinforcement learning strategies, including fine-grained reward shaping and hierarchical policy learning, to further refine reasoning quality. Moreover, we aim to incorporate tool-augmented reasoning, enabling models to dynamically invoke external tools—such as calculators, retrieval systems, or rule engines—during inference, with the goal of improving precision and robustness in high-stakes financial applications.

## Ethics Statement

## Acknowledgements

## References

Jian Chen, Peilin Zhou, Yining Hua, Loh Xin, Kehui Chen, Ziyuan Li, Bing Zhu, and Junwei Liang. FinTextQA: A dataset for long-form financial question answering. In *Proceedings of ACL*, pp. 6025–6047, 2024.

Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. FinQA: A dataset of numerical reasoning over financial data. In *Proceedings of EMNLP*, pp. 3697–3711, 2021.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *CoRR*, abs/2501.12948, 2025. URL https://arxiv.org/abs/2501.12948.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In *Proceedings of NeurIPS*, 2021.

Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, et al. Deepseek-v3 technical report. *CoRR*, abs/2412.19437, 2024. URL https://arxiv.org/abs/2412.19437.

Zhaowei Liu, Xin Guo, Fangqi Lou, Lingfeng Zeng, Jinyi Niu, Zixuan Wang, Jiajie Xu, Weige Cai, Ziwei Yang, Xueqian Zhao, Chao Li, Sheng Xu, Dezhi Chen, Yun Chen, Zuo Bai, and Liwen Zhang. Fin-r1: A large language model for financial reasoning through reinforcement learning. *CoRR*, abs/2503.16252, 2025. URL https://arxiv.org/abs/2503.16252.

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *CoRR*, abs/2501.19393, 2025. URL https://arxiv.org/abs/2501.19393.

OpenAI. Gpt-4 technical report. *CoRR*, abs/2303.08774, 2023. URL https://arxiv.org/abs/2303.08774.

OpenAI. Gpt-4o technical report. https://openai.com/research/gpt-4o, 2024.

OpenAI. Learning to reason with llms. https://openai.com/index/learning-to-reason-with-llms/, 2024.

Lingfei Qian, Weipeng Zhou, Yan Wang, Xueqing Peng, Han Yi, Jimin Huang, Qianqian Xie, and Jianyun Nie. Fino1: On the transferability of reasoning enhanced llms to finance. *CoRR*, abs/2502.08127, 2025. URL https://arxiv.org/abs/2502.08127.

Qwen. QwQ: Reflect Deeply on the Boundaries of the Unknown. https://github.com/QwenLM/QwQ, 2024.

David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a benchmark. In *Proceedings of COLM*, 2024.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *CoRR*, abs/2402.03300, 2024. URL https://arxiv.org/abs/2402.03300.

Qianqian Xie, Weiguang Han, Xiao Zhang, Yanzhao Lai, Min Peng, Alejandro Lopez-Lira, and Jimin Huang. Pixiu: A large language model, instruction data and evaluation benchmark for finance. In *Proceedings of NeurIPS*, pp. 33469–33484, 2023.

Qianqian Xie, Weiguang Han, Zhengyu Chen, Ruoyu Xiang, Xiao Zhang, Yueru He, Mengxi Xiao, Dong Li, Yongfu Dai, Duanyu Feng, Yijing Xu, Haoqiang Kang, Ziyan Kuang, Chenhan Yuan, Kailai Yang, Zheheng Luo, Tianlin Zhang, Zhiwei Liu, Guojun Xiong, Zhiyang Deng, Yuechen Jiang, Zhiyuan Yao, Haohang Li, Yangyang Yu, Gang Hu, Jiajia Huang, Xiao-Yang Liu, Alejandro Lopez-Lira, Benyou Wang, Yanzhao Lai, Hao Wang, Min Peng, Sophia Ananiadou, and Jimin Huang. Finben: A holistic financial benchmark for large language models. In *Proceedings of NeurIPS*, 2024.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, et al. Qwen2.5 technical report. *CoRR*, abs/2412.15115, 2024. URL https://arxiv.org/abs/2412.15115.

Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. Limo: Less is more for reasoning. *CoRR*, abs/2502.03387, 2025. URL https://arxiv.org/abs/2502.03387.

Tianyang Zhong, Zhengliang Liu, Yi Pan, Yutong Zhang, Yifan Zhou, et al. Evaluation of openai o1: Opportunities and challenges of agi. *CoRR*, abs/2409.18486, 2024. URL https://arxiv.org/abs/2409.18486.

Jie Zhu, Junhui Li, Yalong Wen, and Lifan Guo. Benchmarking large language models on CFLUE - a Chinese financial language understanding evaluation dataset. In *Findings of ACL*, pp. 5673–5693, 2024.

## A    Example of converting a multiple-choice question into an open-ended problem

Figure 3 illustrates an example of converting a multiple-choice question into an open-ended question.

## B    Prompts used in this paper

Figure 4 illustrates the prompt used to convert a multiple-choice question from CFLUE into an open-ended format. Since the CFLUE test set includes both single-answer and multiple-answer multiple-choice questions, we design separate prompts for each type. The prompts for single-answer and multiple-answer questions are shown in Figure 5 and Figure 6, respectively.

Figure 7 presents the prompt used to generate answers for FinQA questions.

Figure 8 and Figure 9 show the prompts used to evaluate the correctness of answers for questions in FinQA and CCC, respectively.
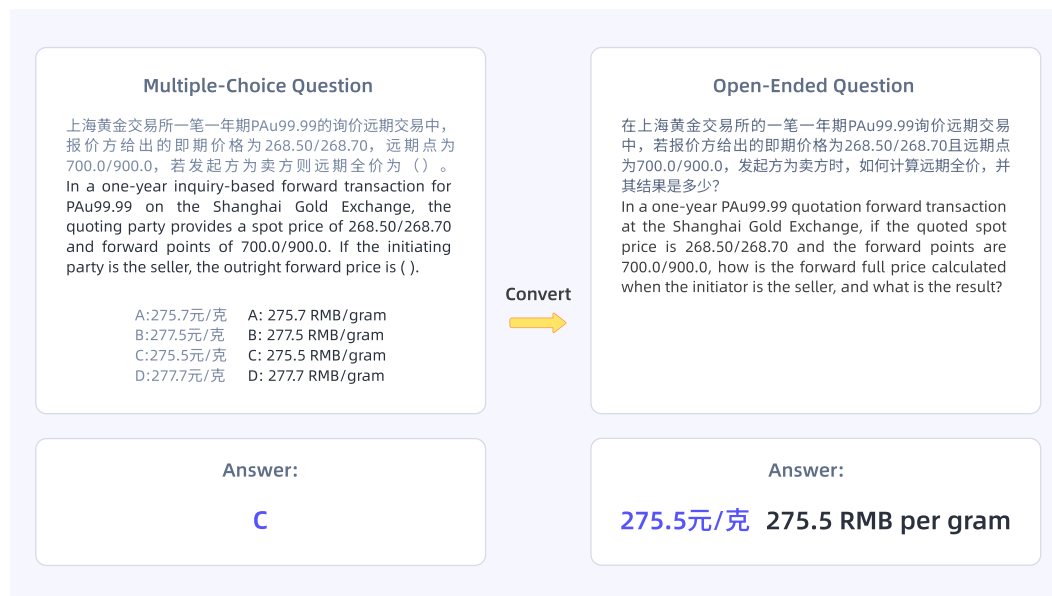
Figure 3: An example of converting a multiple-choice question from CFLUE into an open-ended format.

**Prompt Template**

我将向你提供一个多选题，你的任务是将其改成为一个开放式问题，并提供标准答案。要求如下：

I will provide you with a multiple-choice question, and your task is to rewrite it into an open-ended question and provide a standard answer. The requirements are as follows:

- 问题必须具体，针对原多选题中测试的要点。确保问题是开放式的，即不提供选项，但必须有明确的标准答案。

- The question must be specific and address the key points tested in the original multiple-choice question. Ensure the question is open-ended, meaning no options are provided, but there must be a clear standard answer.

- 根据原问题的正确答案，提供一个简洁的标准答案。答案应当允许精确匹配，以确定模型的回答是否正确。

- Provide a concise standard answer based on the correct answer from the original question. The answer should allow for exact matching to determine if the model's response is correct.

### Multiple-choice Question
{question}
{options}
### Correct Answer
{correct_answer}

## 严格按照以下JSON格式输出：
## Strictly output in the following JSON format:

```
{
    "question": "" # Open-ended Question
    "answer": "" # Standard Answer
}
```

Figure 4: Prompt used to convert a multiple-choice question from CFLUE into an open-ended question.
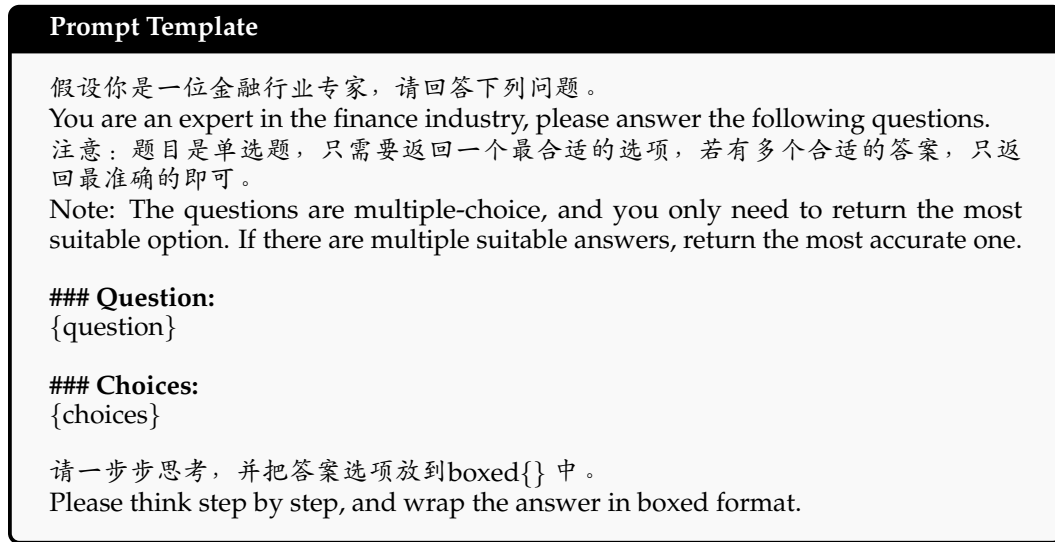
**Prompt Template**

假设你是一位金融行业专家，请回答下列问题。
You are an expert in the finance industry, please answer the following questions.
注意：题目是单选题，只需要返回一个最合适的选项，若有多个合适的答案，只返回最准确的即可。
Note: The questions are multiple-choice, and you only need to return the most suitable option. If there are multiple suitable answers, return the most accurate one.

### Question:
{question}

### Choices:
{choices}

请一步步思考，并把答案选项放到boxed{} 中。
Please think step by step, and wrap the answer in boxed format.

Figure 5: Prompt used to generate answers for single-answer multiple-choice questions in CFLUE.

**Prompt Template**

假设你是一位金融行业专家，请回答下列问题。
You are an expert in the finance industry, please answer the following questions.
注意：题目是多选题，可能存在多个正确的答案。
Note: The questions are multiple-choice and there may be multiple correct answers.

### Question:
{question}

### Choices:
{choices}

请一步步思考，并把答案选项放到boxed{} 中。
Please think step by step, and wrap the answer in boxed format.

Figure 6: Prompt used to generate answers for multiple-answer multiple-choice questions in CFLUE.

**Prompt Template**

Please answer the given financial question based on the context.

### Context:
{context}

### Question:
{question}

Please think step by step, and wrap the answer in boxed{} format.
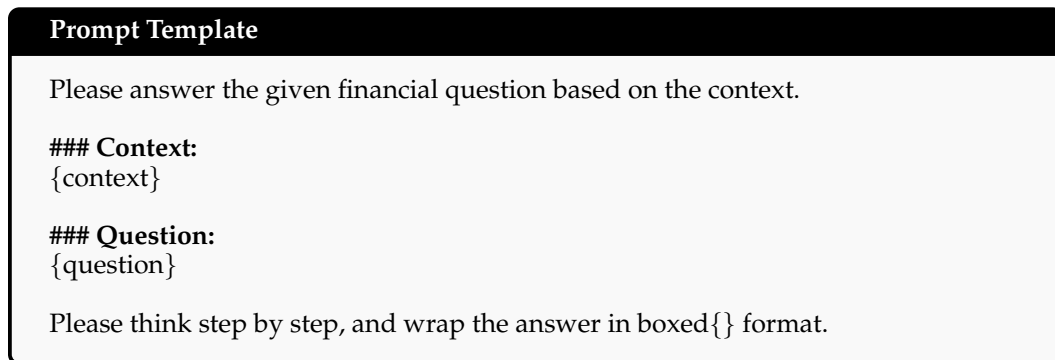
Figure 7: Prompt used to generate answers for questions in FinQA.

**Prompt Template**

You are tasked with comparing a to a to determine if they convey the same meaning based on the following criteria. If they match, output 1; if they do not, output 0.

### Candidate Answer:
{candidate_answer}

### Correct Answer:
{correct_answer}

# Criteria:
- Numerical values are considered consistent if they are the same despite different formats. For example, a of 0.88 is consistent with a of 88%, so return 1.
- Numerical values are also considered consistent if rounding leads to the same result. For example, if the is 8 and the is 7.96, return 1.

# Output:
- Assess the match based on the criteria above and output the result as boxed1 for consistent meanings or boxed0 for inconsistent meanings.

Figure 8: Prompt used to assess the correctness of model-generated answers in FinQA.

**Prompt Template**

# 候选答案(Candidate Answer):
{candidate_answer}

# 正确答案(Correct Answer):
{correct_answer}

# 上述是判断客服是否违规的结论，提供了模型判断的候选答案和相应的正确答案，请根据上述内容，判断候选答案和正确答案是否一致，是否正确。
# The above is a conclusion regarding whether customer service has violated regulations, providing the candidate answer judged by the model and the corresponding correct answer. Please determine whether the candidate answer and the correct answer are consistent and if they are correct.

# 严格按照以下JSON格式输出:
# Strictly output in the following JSON format:

```
{
    "answer": 0/1
}
```

Figure 9: Prompt used to assess the correctness of model-generated answers in CCC.