

What Factors Affect LLMs and RLLMs in Financial Question Answering?

Anonymous ACL submission

Abstract

Recently, large language models (LLMs) and reasoning large language models (RLLMs) have gained considerable attention from many researchers. RLLMs enhance the reasoning capabilities of LLMs through Long Chain-of-Thought (Long CoT) processes, significantly improving the performance of LLMs in addressing complex problems. However, there are few works that systematically explore what methods can fully unlock the performance of LLMs and RLLMs within the financial domain. To investigate the impact of various methods on LLMs and RLLMs, we utilize five LLMs and three RLLMs to assess the effects of prompting methods, agentic frameworks, and multilingual alignment methods on financial question-answering tasks. Our research findings indicate: (1) Current prompting methods and agent frameworks enhance the performance of LLMs in financial question answering by simulating Long CoT; (2) RLLMs possess inherent Long CoT capabilities, which limits the effectiveness of conventional methods in further enhancing their performance; (3) Current advanced multilingual alignment methods primarily improve the multilingual performance of LLMs by extending the reasoning length, which yields minimal benefits for RLLMs. Additionally, we discuss strategies for enhancing the performance of LLMs and RLLMs in financial question answering, which may serve as an inspiration for future improvements. We hope that this study can serve as an important reference for LLMs and RLLMs in the field of financial question answering.

1 Introduction

Recently, large language models (LLMs) have significantly advanced the field of natural language processing (NLP), and more and more researchers utilize LLMs to solve complex tasks in various domains (Zhou et al., 2023; Nie et al., 2024; Qin et al., 2024; Chang et al., 2024; Wang et al., 2025a).

Furthermore, some researchers propose prompting-based and agentic frameworks to enhance the capabilities of LLMs (Wei et al., 2022; Yao et al., 2023a; Zhang et al., 2024b; Yao et al., 2023b; Li et al., 2023; Hong et al., 2024), which further expands the capabilities of LLMs in tasks in different fields. Particularly, to enhance the reasoning capabilities of LLMs in addressing complex problems, researchers introduce reasoning large language models (RLLMs). The application of the Long Chain-of-Thought (Long CoT) fully unlocks the understanding ability of RLLMs (Chen et al., 2025). Long CoT entails deeper reasoning, reflective analysis, and a more extensive exploration of logical structures (Gandhi et al., 2025), which enables RLLMs to perform on par with experts in certain domains.

Researchers have explored the application of LLMs in financial question answering. Wu et al. (2023) present BloombergGPT, a 50 billion parameter LLM for financial domain. Fatemi and Hu (2024) propose a multi-agent framework to enhance the performance of LLMs in financial question answering. Xie et al. (2024) introduce a novel financial benchmark, FinBen, including 42 datasets spanning 24 financial tasks. Srivastava et al. (2024) evaluate LLMs in four financial tabular question-answering datasets. Xue et al. (2024) build a benchmark for financial multilingual multi-modal question answering named FAMMA, which is a challenging benchmark for multi-modal LLMs. Wang et al. (2025c) propose FinSage, which solves the financial filings question answering task through retrieval augmented generation (RAG).

Although LLMs have made significant progress in the field of financial question answering, existing work still mainly focuses on enhancing their performance within this domain. With the advancement of LLMs and RLLMs, the factors influencing their financial question-answering capabilities remain under-explored. Inspired by this gap, we systemati-

085 cally investigate the following issues: (i) the impact
 086 of prompting methods on LLMs/RLLMs; (ii) the
 087 influence of agentic frameworks on LLMs/RLLMs;
 088 and (iii) whether multilingual alignment methods
 089 can enhance the multilingual financial question-
 090 answering abilities of LLMs/RLLMs. We utilize
 091 five LLMs and three RLLMs, conducting detailed
 092 experiments across seven representative methods.

093 The mainly novel insights as follows:

- 094 • **Current prompting methods and agent**
 095 **frameworks enhance the performance of**
 096 **LLMs in financial question answering by**
 097 **simulating Long CoT.** Effective prompting
 098 methods and agent frameworks primarily sim-
 099 ulate Long CoT, enhancing LLM performance
 100 through extended reasoning lengths. This
 101 parallels the performance gains achieved by
 102 RLLMs through Long CoT, demonstrating
 103 that Long CoT represents a significant bot-
 104 tleneck for the current performance improve-
 105 ments of LLMs.
- 106 • **RLLMs possess inherent Long CoT capa-**
 107 **bilities, which limits the effectiveness of**
 108 **conventional methods in further enhancing**
 109 **their performance.** Since RLLMs possess
 110 Long CoT capabilities, the conventional meth-
 111 ods that are effective for LLMs cannot further
 112 enhance RLLM performance. We speculate
 113 that the key to improving RLLMs in the future
 114 lies in the introduction of more complex agent
 115 mechanisms. This would enable RLLMs to
 116 achieve excellent performance through deep
 117 reasoning while also standardizing previous
 118 outputs via the agent framework mechanisms.
- 119 • **Current advanced multilingual alignment**
 120 **methods primarily improve the multilin-**
 121 **gual performance of LLMs by extending**
 122 **the reasoning length, which yields minimal**
 123 **benefits for RLLMs.** Current multilingual
 124 alignment methods depend on the inherent
 125 multilingual capabilities of the model, pri-
 126 marily by extending reasoning lengths to cre-
 127 ate mechanisms similar to Long CoT. This
 128 approach mirrors the gains observed with
 129 prompting methods and agentic frameworks
 130 in LLMs, thus it is not effective for RLLMs.
 131 Consequently, it becomes challenging for
 132 RLLMs to achieve further performance en-
 133 hancements in multilingual contexts.

2 Background 134

135 Financial question answering involves inputting
 136 a question along with relevant context, resulting
 137 in the model generating the answer to the ques-
 138 tion. Specifically, for the multiple-choice question,
 139 the input comprises the question \mathcal{Q} , the available
 140 options $\mathcal{O} = \{l_1, l_2, \dots, l_n\}$, and the context \mathcal{C} ,
 141 while the output consists of the selected options
 142 $\mathcal{A}, \mathcal{A} \in \mathcal{O}$. For the open-ended question, the input
 143 includes the question \mathcal{Q} and the context \mathcal{C} , with the
 144 output being the generated answer \mathcal{A} . These two
 145 types can be defined as follows:

$$146 \mathcal{A}_{choice} = \underset{choice}{\text{Model}}(\mathcal{Q}, \mathcal{O}, \mathcal{C}), \quad (1)$$

$$147 \mathcal{A}_{open} = \underset{open}{\text{Model}}(\mathcal{Q}, \mathcal{C}), \quad (2)$$

148 where $\underset{choice}{\text{Model}}(\cdot)$ presents the model for the
 149 multiple-choice question, $\underset{open}{\text{Model}}(\cdot)$ presents the
 150 model for the open-ended question. 151

3 Exploration 152

153 To thoroughly investigate the impact of various
 154 methods on the performance of LLMs and RLLMs
 155 in financial question answering, this study exam-
 156 ines three perspectives: prompting methods (§ 3.1),
 157 agentic frameworks (§ 3.2), and multilingual align-
 158 ment methods (§ 3.3).

3.1 Exploration of prompting methods 159

160 Current research extensively demonstrates that dif-
 161 ferent prompting methods can significantly alter
 162 the performance of LLMs (Wei et al., 2022; Wang
 163 et al., 2023; Zhang et al., 2024b). The prompting
 164 method modifies certain aspects of the thinking
 165 prompt while ensuring the question remains intact,
 166 enabling LLMs to solve the question by following
 167 the structured steps of the thinking prompt.

168 We select three representative prompting meth-
 169 ods to investigate the impact of varying prompt
 170 word changes on the performance of LLMs and
 171 RLLMs: Direct, Zero-shot CoT (Kojima et al.,
 172 2022), and Plan-and-Solve (Wang et al., 2023).
 173 These prompting methods can be defined as:

- 174 • **Direct:** The direct prompting method enables
 175 LLMs to answer questions by inputting them
 176 directly, without the inclusion of the addi-
 177 tional thinking prompt. This approach allows
 178 for an intuitive assessment of LLM perfor-
 179 mance and minimizes interference from re-
 180 dundant prompts.

181	• Zero-shot CoT: The zero-shot CoT prompting employs the prompt “let’s think step by step” to stimulate the reasoning abilities of LLMs, enabling them to produce longer reasoning processes and, consequently, enhancing their performance.	230
182		231
183		232
184		233
185		234
186		
187	• Plan-and-Solve: The Plan-and-Solve prompting introduces a novel prompt “Let’s first understand the problem and devise a plan to solve the problem. Then, let’s carry out the plan to solve the problem step by step.” This addresses computational errors and enhances the quality of generated reasoning steps through a two-part decomposition.	235
188		236
189		237
190		
191		238
192		239
193		240
194		241
195		242
196		243
197		244
198		245
199		246
200		247
201		248
202		
203		249
204		250
205		251
206		252
207		253
208		254
209		255
210		256
211		257
212		258
213		259
214		260
215		261
216		262
217		263
218		264
219		265
220		266
221		267
222		268
223		269
224		
225		270
226		271
227		272
228		273
229		274

Method	Arithmetic				Non-Arithmetic				Overall
	Overall	Easy	Medium	Hard	Overall	Easy	Medium	Hard	
Llama-3.1-8B-Instruct (Grattafiori et al., 2024)									
Direct (Grattafiori et al., 2024)	10.88	12.66	9.89	9.28	23.44	22.45	32.10	20.27	16.50
Zero-shot CoT (Kojima et al., 2022)	13.20	16.15	10.95	11.07	<u>31.72</u>	32.62	<u>34.73</u>	29.95	21.49
Plan-and-Solve (Wang et al., 2023)	16.09	18.99	<u>14.84</u>	13.17	29.77	32.62	30.52	27.92	22.21
Self-Refine (Madaan et al., 2023)	<u>17.95</u>	<u>22.70</u>	13.07	<u>15.56</u>	32.06	<u>34.74</u>	35.26	<u>29.27</u>	<u>24.26</u>
S ³ Agent (Wang et al., 2024)	19.62	22.27	15.54	19.46	30.80	37.28	32.63	26.57	24.62
GPT-4o-mini (Hurst et al., 2024)									
Direct (Hurst et al., 2024)	32.46	<u>41.70</u>	30.38	21.55	42.52	50.42	50.52	34.90	36.96
Zero-shot CoT (Kojima et al., 2022)	<u>34.13</u>	41.04	34.27	<u>24.55</u>	<u>45.17</u>	<u>56.35</u>	52.10	<u>36.26</u>	<u>39.07</u>
Plan-and-Solve (Wang et al., 2023)	36.18	44.97	<u>33.92</u>	26.04	45.74	56.77	<u>51.05</u>	37.61	40.46
Self-Refine (Madaan et al., 2023)	33.30	41.26	33.21	22.45	42.06	52.11	45.78	35.13	37.22
S ³ Agent (Wang et al., 2024)	30.88	38.42	28.26	22.75	41.72	51.69	46.84	34.23	35.73
Gemini-1.5-flash (Team et al., 2024)									
Direct (Team et al., 2024)	37.20	47.16	35.33	25.14	45.28	<u>58.05</u>	49.47	36.71	40.82
Zero-shot CoT (Kojima et al., 2022)	<u>37.39</u>	<u>46.72</u>	<u>36.04</u>	<u>25.74</u>	<u>46.09</u>	58.89	<u>51.57</u>	<u>36.93</u>	<u>41.28</u>
Plan-and-Solve (Wang et al., 2023)	37.76	46.50	39.57	24.25	47.12	55.50	55.26	39.18	41.95
Self-Refine (Madaan et al., 2023)	34.41	42.35	33.92	23.95	42.52	52.96	50.00	33.78	38.04
S ³ Agent (Wang et al., 2024)	37.20	45.85	34.62	27.54	44.59	57.20	50.52	35.36	40.51
Qwen-2.5-32B (Yang et al., 2024)									
Direct (Yang et al., 2024)	41.86	<u>50.65</u>	40.98	30.53	48.62	59.32	<u>56.31</u>	39.63	44.88
Zero-shot CoT (Kojima et al., 2022)	<u>42.32</u>	50.21	40.98	32.63	50.80	59.32	60.00	<u>42.34</u>	46.11
Plan-and-Solve (Wang et al., 2023)	39.44	44.97	43.46	28.44	<u>49.77</u>	<u>60.59</u>	55.78	41.44	44.06
Self-Refine (Madaan et al., 2023)	41.58	49.56	<u>42.75</u>	29.64	49.65	57.62	54.73	43.24	45.19
S ³ Agent (Wang et al., 2024)	42.51	51.52	41.34	<u>31.13</u>	48.85	61.44	<u>56.31</u>	38.96	<u>45.34</u>
DeepSeek-V3 (Liu et al., 2024)									
Direct (Liu et al., 2024)	56.27	65.50	<u>58.65</u>	41.61	<u>62.06</u>	73.72	69.47	52.70	58.86
Zero-shot CoT (Kojima et al., 2022)	<u>56.18</u>	<u>63.97</u>	59.71	<u>42.51</u>	61.83	<u>71.18</u>	<u>70.52</u>	<u>53.15</u>	58.71
Plan-and-Solve (Wang et al., 2023)	55.16	<u>63.97</u>	56.18	<u>42.21</u>	63.33	70.76	71.05	56.08	<u>58.81</u>
Self-Refine (Madaan et al., 2023)	53.86	61.79	55.83	41.31	59.54	69.06	64.73	52.25	56.40
S ³ Agent (Wang et al., 2024)	54.69	61.57	56.89	43.41	59.42	<u>71.18</u>	66.31	50.22	56.81
DeepSeek-R1-Distill-Qwen-32B (Guo et al., 2025)									
Direct (Guo et al., 2025)	50.41	<u>60.48</u>	50.17	36.82	<u>57.12</u>	65.25	<u>64.73</u>	49.54	53.41
Zero-shot CoT (Kojima et al., 2022)	51.25	60.91	51.94	37.42	56.55	<u>68.64</u>	63.15	47.29	<u>53.62</u>
Plan-and-Solve (Wang et al., 2023)	<u>51.16</u>	60.91	<u>50.53</u>	<u>38.32</u>	55.86	66.52	62.63	47.29	53.26
Self-Refine (Madaan et al., 2023)	44.65	52.18	47.70	31.73	52.06	57.20	60.00	45.94	47.96
S ³ Agent (Wang et al., 2024)	51.25	59.60	51.94	39.22	58.04	69.49	65.26	<u>48.87</u>	54.29
Qwen-3-14B (Yang et al., 2025)									
Direct (Yang et al., 2025)	54.88	66.59	58.65	35.62	<u>59.08</u>	72.03	70.52	47.29	<u>56.76</u>
Zero-shot CoT (Kojima et al., 2022)	55.72	68.34	56.18	38.02	59.19	69.91	70.52	48.64	57.27
Plan-and-Solve (Wang et al., 2023)	<u>54.97</u>	<u>66.81</u>	58.65	35.62	57.81	71.18	<u>68.94</u>	45.94	56.24
Self-Refine (Madaan et al., 2023)	54.41	65.93	56.89	<u>36.52</u>	58.27	<u>71.61</u>	66.84	47.52	56.14
S ³ Agent (Wang et al., 2024)	54.13	65.06	<u>57.24</u>	<u>36.52</u>	58.50	71.18	67.36	<u>47.97</u>	56.09
O4-mini (OpenAI, 2025)									
Direct (OpenAI, 2025)	61.30	70.08	65.01	46.10	70.22	<u>77.54</u>	<u>75.78</u>	63.96	65.29
Zero-shot CoT (Kojima et al., 2022)	62.60	<u>69.86</u>	<u>64.66</u>	50.89	71.37	77.11	77.89	65.54	66.52
Plan-and-Solve (Wang et al., 2023)	61.67	<u>69.86</u>	64.31	48.20	72.41	77.96	77.89	67.11	<u>66.47</u>
Self-Refine (Madaan et al., 2023)	61.20	69.21	60.77	<u>50.59</u>	70.34	75.00	74.73	65.99	65.29
S ³ Agent (Wang et al., 2024)	<u>61.76</u>	69.65	62.19	<u>50.59</u>	<u>71.49</u>	76.69	<u>75.78</u>	<u>66.89</u>	66.11

Table 1: The results of prompting and agentic methods. **Bold number** presents the best result for these methods on the current model. Underline number presents the second-best result for these methods on the current model. **Light yellow** presents the current model is LLM, **Light green** presents the current model is RLLM.

methods on most LLMs. This suggests that employing a well-designed general prompting method effectively enhances the performance of LLMs, even if these methods are not originally proposed for financial question answering.

Larger LLMs exhibit increased robustness to various prompting methods. It is evident that for larger LLMs, such as Gemini-1.5-flash and DeepSeek-V3, the performance gap between different prompting methods is not significant, with performance fluctuations remaining under 2%. This indicates that LLMs with superior performance are less affected by input prompts, demonstrating greater robustness to variations in input.

For RLLMs, competitive results can be achieved without the use of prompting methods. Benefiting from the Long CoT process, RLLMs can generate complex reasoning independently, without deliberate guidance from prompting methods. This renders prompting, which is effective for LLMs, less effective for RLLMs and may even diminish their performance. As shown in Table 1, Plan-and-Solve demonstrates poor performance on the three RLLMs.

4.3 Analysis of Agentic Frameworks

Smaller LLMs can derive greater benefits from the agentic framework. Smaller LLMs, such as

Performance	Direct	Translate-en (Shi et al., 2023)	Self-Refine (Madaan et al., 2023)	S ³ Agent (Wang et al., 2024)	CLP (Qin et al., 2023)
Meta-Llama-3.1-8B-Instruct (Grattafiori et al., 2024)					
Chinese	16.54	16.54 (+0.00)	18.97 (+2.43)	19.22 (+2.68)	23.35 (+6.81)
French	15.38	16.66 (+1.28)	13.46 (-1.92)	21.15 (+5.77)	17.30 (+1.92)
Overall (Chinese+French)	16.22	16.57 (+0.35)	17.46 (+1.24)	19.75 (+3.53)	21.69 (+5.47)
GPT-4o-mini (Hurst et al., 2024)					
Chinese	33.09	34.54 (+1.45)	34.79 (+1.70)	31.63 (-1.46)	38.19 (+5.10)
French	28.84	30.76 (+1.92)	28.84 (+0.00)	32.69 (+3.85)	32.05 (+3.21)
Overall (Chinese+French)	31.92	33.50 (+1.58)	33.15 (+1.23)	31.92 (+0.00)	36.50 (+4.58)
Gemini-1.5-flash (Team et al., 2024)					
Chinese	40.38	38.68 (-1.70)	36.73 (-3.65)	39.65 (-0.73)	37.22 (-3.16)
French	27.56	29.48 (+1.92)	26.92 (-0.64)	27.56 (+0.00)	29.48 (+1.92)
Overall (Chinese+French)	36.86	36.15 (-0.71)	34.03 (-2.83)	36.33 (-0.53)	35.09 (-1.77)
Qwen-2.5-32B (Yang et al., 2024)					
Chinese	41.84	44.52 (+2.68)	44.28 (+2.44)	42.82 (+0.98)	43.30 (+1.46)
French	29.48	29.48 (+0.00)	33.97 (+4.49)	35.89 (+6.41)	35.25 (+5.77)
Overall (Chinese+French)	38.44	40.38 (+1.94)	41.44 (+3.00)	40.91 (+2.47)	41.09 (+2.65)
DeepSeek-V3 (Liu et al., 2024)					
Chinese	57.42	58.63 (+1.21)	55.96 (-1.46)	57.66 (+0.24)	61.31 (+3.89)
French	50.00	50.00 (+0.00)	47.43 (-2.57)	44.23 (-5.77)	57.69 (+7.69)
Overall (Chinese+French)	55.37	56.26 (+0.89)	53.61 (-1.76)	53.96 (-1.41)	60.31 (+4.94)
DeepSeek-R1-Distill-Qwen-32B (Guo et al., 2025)					
Chinese	48.41	49.87 (+1.46)	46.47 (-1.94)	52.06 (+3.25)	50.85 (+2.44)
French	49.35	42.94 (-6.41)	44.87 (-1.16)	49.35 (+0.00)	52.56 (+3.21)
Overall (Chinese+French)	48.67	47.97 (-0.70)	46.03 (-2.64)	51.32 (+2.65)	51.32 (+2.65)
Qwen-3-14B (Yang et al., 2025)					
Chinese	57.66	56.69 (-0.97)	55.23 (-2.43)	55.47 (-2.19)	53.77 (-3.89)
French	48.71	50.00 (+1.29)	42.3 (-6.41)	48.71 (+0.00)	51.28 (+2.57)
Overall (Chinese+French)	55.20	54.85 (-0.35)	51.67 (-3.53)	53.61 (-1.59)	53.08 (-2.12)
O4-mini (OpenAI, 2025)					
Chinese	65.45	65.93 (+0.48)	64.23 (-1.22)	67.63 (+2.18)	64.54 (-0.91)
French	62.17	61.53 (-0.64)	60.25 (-1.92)	61.53 (-0.64)	63.46 (+1.29)
Overall (Chinese+French)	64.55	64.72 (+0.17)	63.13 (-1.42)	65.69 (+1.14)	64.24 (-0.31)

Table 2: The results of multilingual alignment methods. **Bold number** presents the best result for these methods on the current model. **Light yellow** presents the current model is LLM, **Light green** presents the current model is RLLM. The performance of gains/drops relative to the Direct are highlight with **green/red** in the Table.

Llama-3.1-8B-Instruct, can significantly improve their performance through the complex agentic framework, which is shown in Table 1. We speculate that this improvement may stem from the relatively weaker ability of smaller LLMs to understand and follow instructions. The standardized agentic framework helps mitigate the probability of the hallucination, thereby further enhancing overall performance.

A well-designed agentic framework is essential for further enhancing the performance of LLMs. For larger LLMs, directly employing the agentic framework does not lead to significant performance improvements. This may be due to the fact that these agentic frameworks are not specifically designed for financial question answering, and merely transferring them does not fully enhance LLM performance. To enhance performance in financial question answering, a well-designed agentic framework tailored to this task is necessary.

The performance gains of RLLMs within agentic frameworks primarily stem from effective agent collaboration. In RLLMs, agentic frameworks exhibit an improvement compared with prompting methods, particularly in S³ Agent for DeepSeek-R1-Distill-Qwen-32B. This indi-

cates that frameworks successful for LLMs can also be adapted for RLLMs. However, since RLLMs are developed through Long CoT, the performance enhancements associated with agentic frameworks, when employing longer thinking processes, are not pronounced. The benefits of these frameworks are more evident in their inherent design features, such as the reflective capabilities provided by Self-Refine and the multi-faceted thinking fostered by S³ Agent.

4.4 Analysis of Multilingual Alignment Methods

The Translate-en method can enhance multilingual performance for LLMs to some extent; however, the improvement is not substantial. As indicated in Table 2, the Translate-en method yields gains across most models compared to the Direct approach, but these gains remain insignificant. This may be attributed to the fact that Translate-en aligns multilingual tasks to a single language. Consequently, the overall reasoning length for LLMs does not increase, resulting in a lack of depth in thinking, which in turn limits the significant enhancement of performance.

Extending the reasoning length significantly

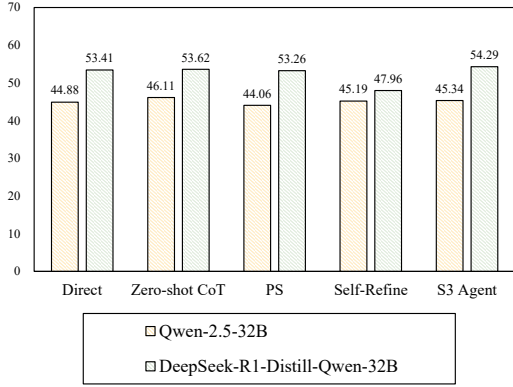


Figure 1: The performance of Qwen-2.5-32B and DeepSeek-R1-Distill-Qwen-32B.

enhances the performance of LLMs. For most LLMs, employing the CLP method yields better performance than the Translate-en approach. Additionally, the agentic framework (Self-Refine and S³ Agent) enhances the performance of LLMs to some degree. This indicates that CLP not only aligns multiple languages to the target language but also further improves performance by strengthening the Long CoT of LLMs. Since the initial stage alignment performance of CLP is also influenced by the inherent multilingual capabilities of the LLMs, insufficient multilingual ability may limit the performance gains achievable through CLP.

RLLMs demonstrate self-alignment capabilities for multilingual questions. In RLLMs, the effectiveness of various alignment methods, such as Translate-en and CLP, is limited. It is posited that RLLMs can achieve results comparable to those of CLP through Long CoT. This suggests that the introduction of CLP may induce overthinking in RLLMs, resulting in a decline in performance. While translation does yield some performance decline, the reduction is not significant.

4.5 The benefits of Long CoT cannot be offset by prompting methods and agentic frameworks

We compare the performance of Qwen-2.5-32B and DeepSeek-R1-Distill-Qwen-32B, where DeepSeek-R1-Distill-Qwen-32B is a Qwen-2.5 model distilled from DeepSeek-R1, as illustrated in Figure 1. It is evident that DeepSeek-R1-Distill-Qwen-32B outperforms the base model across various methods, demonstrating an average improvement of 7.4%. This indicates that the gains from Long CoT significantly enhance the reasoning capabilities of the base model after RLLMs distillation.

Model	Direct	Zero-shot CoT	Plan-and-Solve
Meta-Llama-3.1-8B-Instruct	464.36	973.27	1131.86
gpt-4o-mini	249.77	430.41	531.83
gemini-1.5-flash	286.32	296.03	313.42
Qwen-2.5-32B	368.91	470.99	476.74
DeepSeek-V3	425.65	450.58	493.32
DeepSeek-R1-Distill-Qwen-32B	1,972.60	2006.39	2196.63
Qwen-3-14B	2,087.49	2120.97	2211.74

Table 3: Statistics of average output token consumption in current LLMs/RLLMs.

Consequently, it is challenging for the base model to surpass RLLMs equipped with Long CoT capabilities through advanced frameworks. While it is possible to narrow this gap using a more complex framework, we argue that this approach may lead to token consumption of LLMs approaching that of RLLMs with the same base model.

4.6 Longer reasoning processes generally lead to better performance

To explore the relationship between output length and performance in LLMs and RLLMs, we conduct a statistical analysis of output tokens. The statistics presented in Table 3 indicate that within the same model, more robust methods require a greater number of tokens. This suggests a positive correlation between performance and token usage in the current LLMs, effective prompting techniques enhance performance by simulating a Long CoT. Statistic of Direct in Figure 2 reveals that, in most LLMs, longer reasoning processes with improved performance. This suggests that sufficient cognitive processing benefits LLMs and their current limitations primarily stem from a constrained thinking process.

In contrast, for RLLMs, the output of Long CoT indicates that the length of the CoT is not the primary determinant of performance. For instance, DeepSeek-R1-Distill-Qwen-32B generates an excessive number of tokens for easy questions, yet this does not lead to significant performance enhancements compared with Qwen-3-14B. In easy questions, DeepSeek-R1-Distill-Qwen-32B achieves 53.41%, Qwen-3-14B achieves 56.76% in Acc. This phenomenon may arise from the overthinking in RLLMs (Chen et al., 2024), which fails to yield performance gains while substantially increasing computational costs. Moving forward, dynamically adjusting CoT length according to the complexity of the input problem will represent a critical area of research.

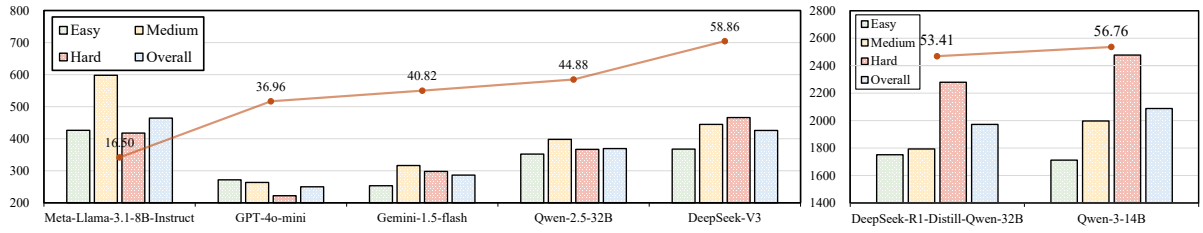
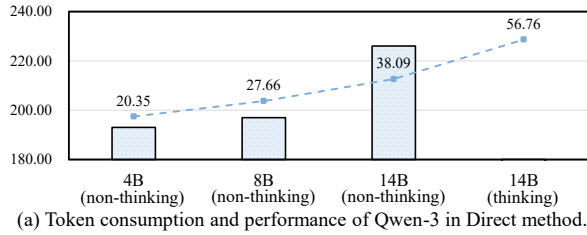
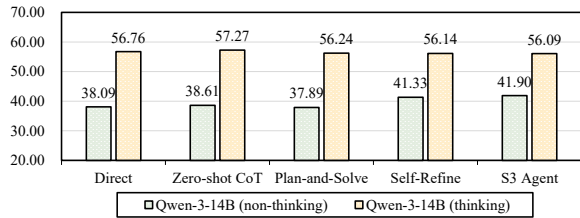


Figure 2: Statistics of average output token consumption for questions of different difficulty in LLMs/RLLMs. The line chart shows the performance of LLMs/RLLMs in the Direct method.



(a) Token consumption and performance of Qwen-3 in Direct method.



(b) Performance of different mode Qwen-3 in different methods.

Figure 3: Performance of different scale and thinking mode in Qwen-3. The line chart in (a) shows Direct performance of Qwen-3, while the histogram in (a) shows average output tokens across scales. The histogram in (b) shows performance in different methods.

4.7 Performance improves as the parameters of LLMs increase

We further analyze the performance and token consumption of Qwen-3 across various scales and the thinking mode, with the final results illustrated in Figure 3. The data reveal that the performance of LLMs adheres to the scaling law (Kaplan et al., 2020), larger models exhibit enhanced performance and longer output lengths. This suggests that more powerful LLMs engage in more extensive reasoning processes, resulting in superior thinking chain capabilities.

Additionally, we examine the performance of the Qwen-3-14B model with thinking mode activated and deactivated. Activating the thinking mode yields a significant improvement, with an average increase of 16.9%. This indicates that RLLMs can achieve superior performance through extended reasoning processes.

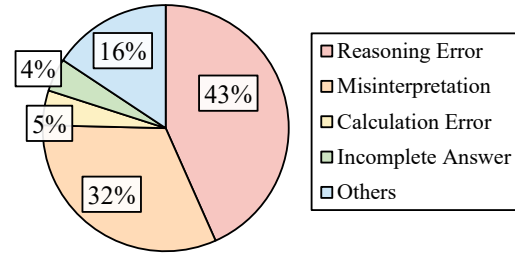


Figure 4: Error statistics for Qwen-3-14B under different thinking modes, considering cases where thinking mode is correct but non-thinking mode fails.

4.8 Understanding the Advantages of Long CoT

To elucidate the advantages conferred by Long CoT, we utilize GPT-4o-mini to systematically analyze the error types observed in Qwen-3-14B under both thinking and non-thinking modes, as illustrated in Figure 4.

The analysis reveals that, in the absence of the thinking mode, the reasoning error constitutes the predominant source of performance degradation. In contrast, activation of the thinking mode substantially mitigates such errors, highlighting the efficacy of Long CoT in enhancing the reasoning capabilities of RLLMs. Furthermore, RLLMs demonstrate superior comprehension of problem statements compared to LLMs, as evidenced by the fact that approximately one-third of LLM errors are attributable to misinterpretation of questions and underlying assumptions. These findings collectively underscore the critical role of Long CoT in advancing the overall performance of financial question answering models.

4.9 It's important to select the optimal alignment language in CLP

We explore the performance differences across various target languages within the CLP method, with the final results presented in Table 4. The find-

Source Language Target Language	English		Chinese		French	
	Chinese	French	English	French	English	Chinese
Meta-Llama-3.1-8B-Instruct	19.88 ▲	14.00 ▼	23.35 ▲	13.38 ▼	17.30 ▲	18.58 ▲
GPT-4o-mini	39.47 ▲	38.53 ▼	38.19 ▲	35.03 ▲	32.05 ▲	28.84 -
Gemini-1.5-flash	41.14 ▼	41.50 ▼	37.22 ▼	38.92 ▼	29.48 ▲	30.76 ▲
Qwen-2.5-32B	45.57 ▼	44.99 ▼	43.30 ▲	42.82 ▲	35.25 ▲	33.97 ▲
DeepSeek-V3	62.55 ▲	60.88 ▲	61.31 ▲	59.36 ▲	57.69 ▲	53.84 ▲
DeepSeek-R1-Distill-Qwen-32B	50.79 ▼	51.45 ▼	50.85 ▲	48.90 ▲	52.56 ▲	51.28 ▲
Qwen-3-14B	57.91 ▲	57.03 ▼	53.77 ▼	58.39 ▲	51.28 ▲	50.00 ▲

Table 4: The results of CLP. ▲ presents the performance of CLP is better than the Direct method, ▼ presents the performance of CLP is worse than the Direct method.

ings indicate that different target languages yield varying performance gains, and directly aligning local languages to English through CLP is often not the most effective strategy. Utilizing various language alignments has demonstrated differing impacts on performance (Wang et al., 2025b), likely attributable to LLMs’ comprehension of distinct language families.

Furthermore, different models utilize distinct optimal target languages identified by the CLP method, which may be attributed to varying language capabilities among the models. Although some researchers propose CLSP (Qin et al., 2023) and AutoCAP (Zhang et al., 2024a). These approaches incur substantial operational costs. Therefore, the selection of the best target language based on the multilingual capabilities of each model requires further exploration.

5 Related Work

With the development of LLMs, more and more researchers focus on the application of LLM in various fields. Improving the performance of LLMs in downstream tasks has consistently been a focal point of researchers’ efforts. Wei et al. (2022) explore the impact of Chain-of-Thought on LLMs and find that forcing thinking of LLMs can significantly improve performance. Inspired by this, researchers further explored the improvement of LLMs. Yao et al. (2023a) propose Tree-of-Thought (ToT), which enhances LLM by considering multiple reasoning paths and self-evaluating. Yin et al. (2023) introduce Exchange-of-Thought (EoT) to enable cross-model communication. In EoT, LLMs could enhance their performance through different network topologies. Zhang et al. (2024b) employs Multi-Perspective Verification and Wrong Information Utilization to prevent LLMs from repeating

the same mistakes, thereby significantly enhancing their performance.

Some researchers have introduced complex agentic frameworks to enhance performance by standardizing the processes of LLMs. Yao et al. (2023b) present a novel agent framework ReAct. ReAct can track the interactions of LLMs with the environment and make decisions about the next steps. Li et al. (2023) propose a novel communicative agent framework that enables LLMs to understand tasks through role-playing. Hong et al. (2024) introduce Standardized Operating Procedures (SOPs) into frameworks and propose MetaGPT, which enhances the ability of software engineering in LLMs. Wang et al. (2024) propose S³ Agent to improve the performance of LLMs on multi-modal problems by introducing multiple perspectives information.

6 Conclusion

In this work, we systematically investigate the factors influencing the performance of LLMs and RLLMs in financial question answering. Through extensive experiments on multiple models and methods, we find that the effectiveness of prompting strategies, agentic frameworks, and multilingual alignment approaches for LLMs largely stems from their ability to simulate longer reasoning chains (Long CoT). RLLMs, equipped with inherent Long CoT capabilities, show limited improvement from these conventional methods, highlighting a performance bottleneck. Our analysis suggests that future advances for RLLMs may require more sophisticated agentic mechanisms and dynamic reasoning processes. We hope these findings deepen our understanding of the Long CoT and serve as a valuable reference in financial question-answering.

549 Limitations

550 Due to computational cost constraints, we primarily
551 selected representative methods for our exper-
552 iments. These methods may not fully reflect the
553 performance of the majority of techniques in finan-
554 cial question answering. To prevent data leakage,
555 we utilized the latest benchmark, FAMMA, for our
556 experiments. However, we cannot guarantee that
557 similar data exists for some questions in the internet
558 data. The challenge of LLMs fitting data requires
559 further exploration in future work.

560 References

561 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,
562 Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi,
563 Cunxiang Wang, Yidong Wang, and 1 others. 2024.
564 A survey on evaluation of large language models.
565 *ACM transactions on intelligent systems and technol-*
566 *ogy*, 15(3):1–45.

567 Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng,
568 Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang
569 Zhou, Te Gao, and Wanxiang Che. 2025. Towards
570 reasoning era: A survey of long chain-of-thought
571 for reasoning large language models. *arXiv preprint*
572 *arXiv:2503.09567*.

573 Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He,
574 Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu,
575 Mengfei Zhou, Zhuosheng Zhang, and 1 others.
576 2024. Do not think that much for $2+3=?$ on
577 the overthinking of o1-like llms. *arXiv preprint*
578 *arXiv:2412.21187*.

579 Sorouralsadat Fatemi and Yuheng Hu. 2024. [Enhancing](#)
580 [financial question answering with a multi-agent re-](#)
581 [flection framework](#). In *Proceedings of the 5th ACM*
582 *International Conference on AI in Finance, ICAIF*
583 *'24*, pages 530–537, New York, NY, USA. Associa-
584 tion for Computing Machinery.

585 Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh,
586 Nathan Lile, and Noah D Goodman. 2025. Cognitive
587 behaviors that enable self-improving reasoners, or,
588 four habits of highly effective stars. *arXiv preprint*
589 *arXiv:2503.01307*.

590 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
591 Abhinav Pandey, Abhishek Kadian, Ahmad Al-
592 Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,
593 Alex Vaughan, and 1 others. 2024. The llama 3 herd
594 of models. *arXiv preprint arXiv:2407.21783*.

595 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao
596 Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-
597 rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.
598 Deepseek-r1: Incentivizing reasoning capability in
599 llms via reinforcement learning. *arXiv preprint*
600 *arXiv:2501.12948*.

Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiauwu
Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang,
Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang
Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu,
and Jürgen Schmidhuber. 2024. [MetaGPT: Meta pro-](#)
[gramming for a multi-agent collaborative framework](#).
In *The Twelfth International Conference on Learning*
Representations.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam
Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow,
Akila Welihinda, Alan Hayes, Alec Radford, and 1
others. 2024. Gpt-4o system card. *arXiv preprint*
arXiv:2410.21276.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B
Brown, Benjamin Chess, Rewon Child, Scott Gray,
Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.
Scaling laws for neural language models. *arXiv*
preprint arXiv:2001.08361.

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

Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii
Khizbullin, and Bernard Ghanem. 2023. Camel:
Communicative agents for "mind" exploration of
large language model society. *Advances in Neural*
Information Processing Systems, 36:51991–52008.

Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,
Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi
Deng, Chenyu Zhang, Chong Ruan, and 1 others.
2024. Deepseek-v3 technical report. *arXiv preprint*
arXiv:2412.19437.

Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao,
Junwei Yang, Yiyang Gu, Bohan Wu, Binqi Chen,
Ziyue Qiao, Qingqing Long, and 1 others. 2025a.
Large language model agent: A survey on method-
ology, applications and challenges. *arXiv preprint*
arXiv:2503.21460.

Xufang Luo, Yuge Zhang, Zhiyuan He, Zilong Wang,
Siyun Zhao, Dongsheng Li, Luna K Qiu, and
Yuqing Yang. 2025b. Agent lightning: Train any
ai agents with reinforcement learning. *arXiv preprint*
arXiv:2508.03680.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler
Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,
Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,
and 1 others. 2023. Self-refine: Iterative refinement
with self-feedback. *Advances in Neural Information*
Processing Systems, 36:46534–46594.

Yuqi Nie, Yaxuan Kong, Xiaowen Dong, John M Mul-
vey, H Vincent Poor, Qingsong Wen, and Stefan
Zohren. 2024. A survey of large language models
for financial applications: Progress, prospects and
challenges. *arXiv preprint arXiv:2406.11903*.

656	OpenAI. 2025. Openai o3 and o4-mini system card . Accessed: 2025-07-07.	Xinyu Wang, Jijun Chi, Zhenghan Tai, Tung Sum Thomas Kwok, Muzhi Li, Zhuhong Li, Hailin He, Yuchen Hua, Peng Lu, Suyuchen Wang, and 1 others. 2025c. Finsage: A multi-aspect rag system for financial filings question answering. <i>arXiv preprint arXiv:2504.14493</i> .	711
657			712
658	Libo Qin, Qiguang Chen, Xiachong Feng, Yang Wu, Yongheng Zhang, Yinghui Li, Min Li, Wanxiang Che, and Philip S Yu. 2024. Large language models meet nlp: A survey. <i>arXiv preprint arXiv:2405.12819</i> .		713
659			714
660			715
661			716
662	Libo Qin, Qiguang Chen, Fuxuan Wei, Shijue Huang, and Wanxiang Che. 2023. Cross-lingual prompting: Improving zero-shot chain-of-thought reasoning across languages . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2695–2709, Singapore. Association for Computational Linguistics.	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	717
663			718
664			719
665			720
666			721
667			722
668		Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhajan Kambadur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. <i>arXiv preprint arXiv:2303.17564</i> .	723
669	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. 2023. Language models are multi-lingual chain-of-thought reasoners . In <i>The Eleventh International Conference on Learning Representations</i> .		724
670			725
671			726
672			727
673		Qianqian Xie, Weiguang Han, Zhengyu Chen, Ruoyu Xiang, Xiao Zhang, Yueru He, Mengxi Xiao, Dong Li, Yongfu Dai, Duanyu Feng, and 1 others. 2024. Finben: A holistic financial benchmark for large language models. <i>Advances in Neural Information Processing Systems</i> , 37:95716–95743.	728
674			729
675			730
676	Pragya Srivastava, Manuj Malik, Vivek Gupta, Tanuja Ganu, and Dan Roth. 2024. Evaluating llms’ mathematical reasoning in financial document question answering. <i>arXiv preprint arXiv:2402.11194</i> .		731
677			732
678			733
679		Siqiao Xue, Tingting Chen, Fan Zhou, Qingyang Dai, Zhixuan Chu, and Hongyuan Mei. 2024. Famma: A benchmark for financial domain multilingual multimodal question answering. <i>arXiv preprint arXiv:2410.04526</i> .	734
680	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .		735
681			736
682			737
683			738
684		An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	739
685			740
686	Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2609–2634.		741
687			742
688			743
689			744
690			745
691		An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024. Qwen2.5 technical report. <i>arXiv preprint arXiv:2412.15115</i> .	746
692			747
693	Peng Wang, Wenpeng Lu, Chunlin Lu, Ruoxi Zhou, Min Li, and Libo Qin. 2025a. Large language model for medical images: A survey of taxonomy, systematic review, and future trends . <i>Big Data Mining and Analytics</i> , 8(2):496–517.		748
694			749
695		Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023a. Tree of thoughts: Deliberate problem solving with large language models. <i>Advances in neural information processing systems</i> , 36:11809–11822.	750
696			751
697			752
698	Peng Wang, Ruihan Tao, Qiguang Chen, Mengkang Hu, and Libo Qin. 2025b. X-WebAgentBench: A multilingual interactive web benchmark for evaluating global agentic system . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 19320–19335, Vienna, Austria. Association for Computational Linguistics.		753
699			754
700		Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023b. React: Synergizing reasoning and acting in language models. In <i>International Conference on Learning Representations (ICLR)</i> .	755
701			756
702			757
703			758
704			759
705	Peng Wang, Yongheng Zhang, Hao Fei, Qiguang Chen, Yukai Wang, Jiasheng Si, Wenpeng Lu, Min Li, and Libo Qin. 2024. S3 agent: Unlocking the power of vllm for zero-shot multi-modal sarcasm detection. <i>ACM Trans. Multimedia Comput. Commun. Appl.</i> Just Accepted.	Zhangyue Yin, Qiushi Sun, Cheng Chang, Qipeng Guo, Junqi Dai, Xuanjing Huang, and Xipeng Qiu. 2023. Exchange-of-thought: Enhancing large language model capabilities through cross-model communication . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 15135–15153, Singapore. Association for Computational Linguistics.	760
706			761
707			762
708			763
709			764
710			765
			766
			767

768	Guibin Zhang, Hejia Geng, Xiaohang Yu, Zhenfei Yin,	processes (Luo et al., 2025b; Zhang et al., 2025).	821
769	Zaibin Zhang, Zelin Tan, Heng Zhou, Zhongzhi Li,	Developing a suitable agentic workflow and further	822
770	Xiangyuan Xue, Yijiang Li, and 1 others. 2025. The	training LLMs and RLLMs through agentic rein-	823
771	landscape of agentic reinforcement learning for llms:	forcement learning represents an intriguing avenue	824
772	A survey. <i>arXiv preprint arXiv:2509.02547</i> .	for improving model performance.	825
773	Yongheng Zhang, Qiguang Chen, Min Li, Wanxiang		
774	Che, and Libo Qin. 2024a. Autocap: Towards auto-		
775	matic cross-lingual alignment planning for zero-shot		
776	chain-of-thought. In <i>Findings of the Association for</i>		
777	<i>Computational Linguistics ACL 2024</i> , pages 9191–		
778	9200.		
779	Yongheng Zhang, Qiguang Chen, Jingxuan Zhou, Peng		
780	Wang, Jiasheng Si, Jin Wang, Wenpeng Lu, and Libo		
781	Qin. 2024b. Wrong-of-thought: An integrated rea-		
782	soning framework with multi-perspective verification		
783	and wrong information . In <i>Findings of the Associa-</i>		
784	<i>tion for Computational Linguistics: EMNLP 2024</i> ,		
785	pages 6644–6653, Miami, Florida, USA. Association		
786	for Computational Linguistics.		
787	Hongjian Zhou, Fenglin Liu, Boyang Gu, Xinyu		
788	Zou, Jinfa Huang, Jinge Wu, Yiru Li, Sam S		
789	Chen, Peilin Zhou, Junling Liu, and 1 others. 2023.		
790	A survey of large language models in medicine:		
791	Progress, application, and challenge. <i>arXiv preprint</i>		
792	<i>arXiv:2311.05112</i> .		
793	A The Use of Large Language Models		
794	(LLMs)		
795	We declare that only LLMs were utilized to polish		
796	the English of this paper.		
797	B Discussion		
798	While considerable research has focused on enhanc-		
799	ing the performance of LLMs across various do-		
800	main, our experiments indicate that the key to im-		
801	proving LLM performance may lie in strengthening		
802	their long-term thinking capabilities, which aligns		
803	with the overarching objective of RLLMs. Conse-		
804	quently, we find that currently effective prompting		
805	methods and agentic frameworks for LLMs do not		
806	yield significant improvements for RLLMs. It is		
807	important to note that the agentic framework em-		
808	ployed in our study was not specifically tailored		
809	for financial question answering. Thus, agentic		
810	frameworks that integrate RLLMs with techniques		
811	such as retrieval-augmented generation (RAG) and		
812	are well-designed for financial question answering		
813	may present a promising solution.		
814	Furthermore, the Long CoT approach constrains		
815	the potential for enhancing performance of RLLMs		
816	through agentic frameworks. The advent of agen-		
817	tic reinforcement learning (Agentic RL) enables		
818	LLMs and RLLMs to function not only as stan-		
819	dalone foundation models but also to be incorpo-		
820	rated into complex reasoning and decision-making		