

Can GPT models be Financial Analysts?

An Evaluation of ChatGPT and GPT-4 on mock CFA Exams

Ethan Callanan^{1,†}, Amarachi Mbakwe^{2,†,‡}, Antony Papadimitriou^{3,†}, Yulong Pei^{3,†}, Mathieu Sibue^{3,†}, Xiaodan Zhu¹, Zhiqiang Ma³, Xiaomo Liu³, and Sameena Shah³

¹Queen’s University

²Virginia Tech

³J.P. Morgan AI Research

¹{e.callanan,xiaodan.zhu}@queensu.ca, ²bmamarachi@vt.edu, ³{first.last}@jpmchase.com

Abstract

Large Language Models (LLMs) have demonstrated remarkable performance on a wide range of Natural Language Processing (NLP) tasks, often matching or even beating state-of-the-art task-specific models. This study aims at assessing the financial reasoning capabilities of LLMs. We leverage mock exam questions of the Chartered Financial Analyst (CFA) Program to conduct a comprehensive evaluation of ChatGPT¹ and GPT-4² in financial analysis, considering Zero-Shot (ZS), Chain-of-Thought (CoT), and Few-Shot (FS) scenarios. We present an in-depth analysis of the models’ performance and limitations, and estimate whether they would have a chance at passing the CFA exams. Finally, we outline insights into potential strategies and improvements to enhance the applicability of LLMs in finance. In this perspective, we hope this work paves the way for future studies to continue enhancing LLMs for financial reasoning through rigorous evaluation.

1 Introduction

NLP has undergone a profound transformation driven by the emergence of LLMs. Models such as LLaMA (Touvron et al., 2023), PaLM (Chowdhery et al., 2022), ChatGPT and GPT-4 from OpenAI, have garnered significant attention from both the research community and the public thanks to their impressive dialog style. From text summarization (Barros et al., 2023; Zhang et al., 2023) and code generation (Le et al., 2022; Yang et al., 2023a) to question answering (Wang et al., 2022; Khashabi et al., 2020), named entity recognition (Li et al., 2022), and beyond, LLMs have showcased remarkable performance across diverse tasks.

¹<https://platform.openai.com/docs/models/gpt-3-5>

²<https://platform.openai.com/docs/models/gpt-4>

[†]Equal contribution

[‡]Work done while interning at J.P. Morgan AI Research

Model	Setting	Level I	Level II
ChatGPT	ZS	58.8 ± 0.2	46.6 ± 0.6
	CoT	58.0 ± 0.2	47.2 ± 0.3
	2S	63.0 ± 0.2	46.6 ± 0.1
	4S	62.3 ± 0.2	45.7 ± 0.2
	6S	62.2 ± 0.2	47.0 ± 0.3
	10S	62.4 ± 0.2	47.6 ± 0.4
GPT-4	ZS	73.2 ± 0.2	57.4 ± 1.5
	CoT	74.0 ± 0.2	61.4 ± 0.9
	2S	73.9 ± 0.1	60.2 ± 0.9
	4S	73.8 ± 0.2	60.5 ± 0.7
	6S	74.5 ± 0.2	-
	10S	74.6 ± 0.2	-

Table 1: Overall Performance of ChatGPT and GPT-4 on Level I and Level II Exams (Accuracy) in Zero-Shot (ZS), Chain-of-Thought (CoT), and Few-Shot (FS) settings.

In finance, NLP has played a pivotal role in enhancing various services, such as customer relations, stock sentiment analysis, financial question answering (Wang et al., 2022), document understanding (Kim et al., 2022), and report summarization (Abdaljalil and Bouamor, 2021). Despite these advancements, applying NLP in finance poses unique challenges, such as the distinct nature of financial tasks, linguistic structures, and specialized terminology. As a result, the performance of general NLP models often falls short when applied to finance-related tasks – the specific challenges of financial reasoning problems warrant further investigation.

In this paper, we rigorously assess the capabilities of LLMs in real-world financial reasoning problems by conducting an evaluation on mock exam questions of the prestigious Chartered Financial Analyst (CFA) Program³. The CFA exams

³<https://www.cfainstitute.org/en/programs/cfa/exam>

Which of the following is most likely an assumption of technical analysis?

- A. Security markets are efficient*
- B. Market trends reflect irrational human behavior*
- C. Equity markets react quickly to inflection points in the broad economy*

(a) Level I sample question

Paris Rousseau, a wealth manager at a US-based investment management firm, is meeting with a new client. The client has asked Rousseau to make recommendations regarding his portfolio's exposure to liquid alternative investments [...]
[Table Evidence]

The AFFO per share for Autier REIT over the last 12-months is closest to:

- A. \$6.80.*
- B. \$7.16.*
- C. \$8.43.*

(b) Level II sample question

Figure 1: CFA example questions (source: CFA Institute); the question appears in bold, the multiple choices in blue and italic, and the vignette/case description in orange and italic

are known for their meticulous yet practical assessment of financial expertise, making their resolution an ideal use case to gauge the capabilities of LLMs in handling complex financial reasoning scenarios. Our work focuses on two closed-source, non-domain specific LLMs, ChatGPT and GPT-4, using various popular prompting techniques. Our contributions are as follows:

- 1 We conduct the first comprehensive evaluation of ChatGPT and GPT-4 in financial reasoning problems using CFA mock exam questions, considering ZS, CoT, and FS scenarios.
- 2 We present an in-depth analysis of the models' performance and limitations in solving these financial reasoning problems, and estimate how they would fare in the Level I and Level II CFA exams.
- 3 We outline insights into potential strategies and improvements to enhance the applicability of LLMs in finance, opening new avenues for research and development.

2 Related Work

2.1 LLMs and Finance

LLMs are Transformer-based generative models (Vaswani et al., 2017) trained on massive datasets that cover a broad range of topics and domains. Previous work has demonstrated the ability of LLMs to generalize surprisingly well to unseen downstream tasks, with little to no additional training data (Brown et al., 2020; Wei et al., 2022).

This raises an interesting question on the competitiveness of LLMs with supervised state-of-the-art models on specialized domains, such as finance. Indeed, the characteristics of most financial tasks — which rely on very specific concepts and mathematical formula, frequently leverage diagrams and tables, often need multistep reasoning with calculations — make finance a challenging domain of application for LLMs. Several paths have been proposed to incorporate or emphasize domain-specific knowledge in LLMs: continued pre-training (Araci, 2019; Wu et al., 2023) and supervised fine-tuning on new data (Mosbach et al., 2023; Yang et al., 2023b), retrieval augmented generation using a vector database of external knowledge (Lewis et al., 2020), etc. However, before considering such enhancements, only few papers have proceeded to extensively benchmark the out-of-the-box capabilities of newer instruction-tuned LLMs *in finance* (Li et al., 2023).

2.2 Evaluation of LLMs on Human Exams and other Benchmarks

Several previous studies have evaluated the performance of LLMs on different standard exams. Tests considered include the United States medical licensing exam (Kung et al., 2023), free-response clinical reasoning exams (Strong et al., 2023), college-level scientific exams (Wang et al., 2023), the Bar exam (Katz et al., 2023), the driver's license knowledge test (Rahimi et al., 2023), and more. The crucial contribution that these works bring to the scientific community and the industry is an in-depth analysis of the strengths and weaknesses of LLMs in realistic domain-specific settings. Through their conclusions, such investiga-

Code available at <https://github.com/e-cal/gpt-cfa>

Topic	Level I			Level II		
	Calculations	#Tables	Len(Prompt)	Calculations	#Tables	Len(Prompt)
Ethics	0.7%	0.01	125	0.0%	0.00	1013
Derivatives	20.7%	0.00	65	75.0%	2.00	816
Alternative Investments	36.4%	0.06	85	66.7%	2.00	840
Portfolio Management	38.3%	0.18	110	56.3%	2.13	1077
Fixed Income	43.0%	0.06	87	50.0%	1.45	779
Economics	50.6%	0.25	121	66.7%	2.00	1115
Equity	52.5%	0.19	112	45.8%	1.00	1053
Corporate Issuers	59.3%	0.28	120	44.4%	1.67	930
Quantitative Methods	70.5%	0.26	131	27.8%	0.00	1256
Financial Reporting	57.7%	0.35	151	53.6%	2.79	1383
Overall	42.4%	0.17	116	45.5%	1.47	1058

Table 2: Question characteristics by topic; percentage of questions requiring calculation, average number of table evidence per question, and average prompt length (estimated using the tiktoken Python package)

tions guide subsequent research and practical use case resolutions in industry.

For example, (Wang et al., 2023) evaluated ChatGPT and GPT-4 on a collection of Physics, Chemistry, and Math problems, and then concluded that current LLMs do not deliver satisfactory performance in complex scientific reasoning yet to be reliably leveraged in practice. In contrast, (Bang et al., 2023) found that ChatGPT outperformed fine-tuned task-specific models on four different NLP tasks, thus suggesting ChatGPT could be directly applied to solve industry use cases involving these tasks.

Our paper aims at following the footsteps of (Li et al., 2023) and delves further into the assessment of the inner financial reasoning abilities of ChatGPT and GPT-4 to help future industry applications.

3 Dataset

The CFA Program is a three-part exam that tests the fundamentals of investment tools, valuing assets, portfolio management, and wealth planning. It is typically completed by those who want to work in the financial industry with backgrounds in finance, accounting, economics, or business. Successfully completing the CFA Program reflects a strong grasp of fundamental financial knowledge, and charterholders are then qualified for roles related to investment management, risk management, asset management, and more.

As mentioned above, the CFA exam is composed of three levels, each with a specific format. Irrespective of the level, each problem from the CFA exam is affiliated to one of ten distinct finance topics: Ethics, Quantitative Methods,

Economics, Financial Statement Analysis, Corporate Issuers, Portfolio Management, Equity Investments, Fixed Income, Derivatives, and Alternative Investments. Level I features a total of 180 independent Multiple Choice Questions (MCQs). Level II consists of 22 item sets comprised of vignettes (i.e., case descriptions with evidence) and 88 accompanying MCQs. Finally, Level III comprises a mix of vignette-supported essay questions and vignette-supported multiple choice questions.

Two main challenges arise when trying to benchmark any model on the CFA exam. Firstly, the CFA Institute refrains from publicly releasing past exams taken by registered candidates, making the collection of official questions and answers directly from any CFA exam impossible. Secondly, a significant fraction of the level III item sets expects plain text responses, which then require the costly intervention of human experts for grading. To circumvent these difficulties, we decide to rely on mock CFA exams and choose to solely focus on levels I and II, leaving Level III to future work. We collected a total of five Level I mock exams and two Level II mock exams. We share in Figure 1 example MCQs published by the CFA Institute for Level I and Level II. We ensure each topic is represented in similar proportions to the original CFA sections (Figure 2 and Figure 3 in the Appendix). Table 2 summarizes important statistics about Level I and Level II problems.

4 Experiments

4.1 Setup

This section outlines the methodology employed to assess the financial reasoning abilities of

Category	ChatGPT			GPT-4		
	ZS	CoT	2S	ZS	CoT	10S
Ethics	59.2 ± 0.1	59.2 ± 1.4	64.6 ± 0.9	80.3 ± 0.7	78.9 ± 0.4	82.4 ± 0.5
Quantitative Methods	53.9 ± 0.2	50.0 ± 0.8	59.7 ± 1.0	78.0 ± 0.7	76.0 ± 1.1	76.0 ± 0.8
Economics	68.0 ± 1.1	63.7 ± 2.5	68.0 ± 3.9	74.1 ± 1.9	73.6 ± 1.2	76.2 ± 0.6
Financial Reporting	54.0 ± 1.2	53.4 ± 0.6	60.1 ± 0.7	68.2 ± 1.0	70.8 ± 1.3	70.0 ± 0.7
Corporate Issuers	71.4 ± 5.2	69.8 ± 4.8	74.2 ± 4.1	74.4 ± 4.1	74.6 ± 6.2	75.3 ± 4.0
Equity Investments	59.4 ± 0.1	60.9 ± 0.7	62.5 ± 1.0	80.3 ± 0.7	70.5 ± 0.9	68.8 ± 0.8
Fixed Income	55.6 ± 1.4	60.2 ± 0.5	63.6 ± 0.5	74.9 ± 2.6	60.2 ± 0.5	73.6 ± 0.8
Derivatives	61.1 ± 4.1	68.5 ± 2.1	73.0 ± 1.5	90.5 ± 0.8	93.8 ± 0.7	96.0 ± 0.5
Alternative Investments	60.7 ± 2.4	60.7 ± 1.9	62.9 ± 1.1	75.9 ± 1.1	77.1 ± 1.0	72.1 ± 1.3
Portfolio Management	58.3 ± 2.8	48.3 ± 3.6	61.7 ± 2.4	63.7 ± 0.6	71.7 ± 0.9	79.6 ± 1.4
Overall	58.8 ± 0.2	58.0 ± 0.2	63.0 ± 0.2	73.2 ± 0.2	74.0 ± 0.9	74.6 ± 0.2

Table 3: ChatGPT and GPT-4 accuracy on Level I Exams

Category	ChatGPT			GPT-4		
	ZS	CoT	10S	ZS	CoT	4S
Ethics	31.3 ± 7.6	37.5 ± 9.5	21.9 ± 4.6	43.8 ± 1.6	56.3 ± 1.2	59.4 ± 1.5
Quantitative Methods	44.4 ± 12.0	55.6 ± 6.5	54.2 ± 9.3	66.7 ± 1.1	66.7 ± 7.4	72.2 ± 4.3
Economics	66.7 ± 0.0	58.3 ± 1.4	62.5 ± 1.9	41.7 ± 1.4	58.3 ± 6.3	50.0 ± 6.9
Financial Reporting	39.6 ± 3.4	31.3 ± 2.0	44.8 ± 2.5	54.2 ± 3.9	66.7 ± 4.2	63.5 ± 3.3
Corporate Issuers	55.6 ± 3.7	50.0 ± 2.8	50.0 ± 1.9	77.8 ± 0.9	77.8 ± 0.6	80.6 ± 1.3
Equity Investments	60.4 ± 1.6	60.4 ± 9.9	60.9 ± 7.0	65.0 ± 5.7	58.8 ± 7.3	62.5 ± 4.7
Fixed Income	38.9 ± 0.9	27.8 ± 6.5	34.4 ± 1.9	60.0 ± 5.8	62.2 ± 0.8	53.9 ± 1.9
Derivatives	50.0 ± 5.6	58.3 ± 12.5	47.9 ± 3.1	66.7 ± 5.6	58.3 ± 0.7	50.0 ± 4.2
Alternative Investments	33.3 ± 0.0	33.3 ± 0.0	58.3 ± 0.7	66.7 ± 0.0	50.0 ± 0.0	75.0 ± 0.7
Portfolio Management	47.2 ± 0.9	66.7 ± 8.3	59.7 ± 9.5	36.1 ± 1.6	55.6 ± 0.6	56.9 ± 4.3
Overall	46.6 ± 0.6	47.2 ± 0.3	47.6 ± 0.4	57.4 ± 1.5	61.4 ± 0.9	60.5 ± 0.7

Table 4: ChatGPT and GPT-4 accuracy on Level II Exams

ChatGPT and GPT-4 using mock CFA exams. Our study examined various prompting paradigms.

ZS prompting: We gauged the models’ inherent reasoning abilities without providing any correct examples in the input.

FS prompting: We furnished the models with prior examples of expected behavior to facilitate the acquisition of new knowledge that could aid in solving the questions. We tested two different strategies to select FS examples: (a) randomly sampling from the entire set of questions within the exam level (2S, 4S and 6S), and (b) sampling one question from each topic in the exam level (10S). This last approach aims at enabling the models to discern the distinct attributes of each topic within every exam level. Due to the limited context window of GPT-4 and the length of the Level II item-sets (case description and question), 6S and 10S prompting were not evaluated for GPT-4 on the Level II mock exams.

CoT prompting: For each exam level, we also

evaluated the models by prompting them to think through the input problem step-by-step and show their work for calculations (also known as ZS CoT) (Wei et al., 2022). This has the added benefit of allowing us to analyze the models’ "problem-solving process" and thus determine where and why it might have gone wrong.

Implementation Details: We conducted the experiments using the OpenAI ChatCompletion API (gpt-3.5-turbo and gpt-4 models) with functions and set the temperature parameter to zero, thereby eliminating randomness in the models’ generations. The prompt templates we crafted for each level and for each prompting setting can be found in the Appendix. We employed a memorization test as in (Kıçman et al., 2023) to confirm that the models had not memorized the mock exams as part of their training data.

Metrics: To measure the performance of LLMs on the mock exam MCQs, we compared their predictions against the established solution set of each of the CFA mock exams collected. Accuracy served

as our sole evaluation metric throughout this study.

4.2 Results Overview

LLMs struggle more on Level II than on Level I: We notice that, no matter the prompting paradigm employed, both `ChatGPT` and `GPT-4` encounter more difficulties correctly answering the item-sets from Level II than the independent questions from Level I (Table 3, Table 4). While there is no general consensus as to which level is usually considered harder for exam takers, we suggest that three factors might have negatively affected the performance of LLMs in Level II based on our analysis.

Firstly, the case description attached to each item-set from Level II increases the length of the input prompt and dilutes the useful information it contains. Indeed, we observe that Level II prompts are on average $\sim 10\times$ longer than Level I prompts; confronting Table 2, Table 3, Table 4 shows that topics associated with poor performance usually present longer contexts both in Level I and Level II. In addition, the detailed case descriptions from Level II depict realistic day-to-day situations that contrast with the more general questions from Level I: LLMs thus need to abstract from case-specific details in Level II questions so as to identify the underlying finance concepts involved.

Secondly, as Level II questions are grouped into item-sets, each item tends to go more in-depth about a specific finance topic than the questions that compose Level I, thus leading to more specialized and intricate problems.

Lastly, the Level II section features a slightly higher proportion of questions requiring calculations and a much higher proportion of questions containing table evidence, in comparison to Level I (Table 2). Given the known limitations of LLMs for out-of-the-box numerical and table reasoning (Frieder et al., 2023; Chen et al., 2022), this could also explain the lower accuracy observed in Level II across the board.

GPT-4 outperforms ChatGPT in almost all experiments, but certain finance topics remain challenging for both models: As shown in Table 3 and Table 4, `GPT-4` consistently beats `ChatGPT` in all topics in Level I and most topics in Level II, irrespective of the learning paradigm.

In Level I, we see that both LLMs perform best in the Derivatives, Alternative Investments, Cor-

porate Issuers, Equity Investments, and Ethics topics. For Derivatives and Ethics, this observation can be explained by the low amount of calculations and table understanding required to answer correctly (Table 2). The explicit mention of popular finance notions in the questions of Derivatives and Ethics (e.g., options, arbitrage, etc.) further reduces their difficulty too. Similarly, in Alternative Investments, Corporate Issuers, and Equity Investments, problems often directly refer to well-known finance concepts that might have been encountered by `ChatGPT` and `GPT-4` during pretraining or instruction-tuning – thus facilitating their resolution despite having more calculations involved. However, both models show relatively poor performance in the Financial Reporting and Portfolio Management topics in Level I, with `ChatGPT` also struggling a lot more on highly computational topics such as Quantitative Methods. Indeed, Portfolio Management and Financial Reporting problems are more case-based, applied, computational, and CFA-specific than the ones from the aforementioned topics, which might have negatively affected performance. They also tend to include more table evidence and complex details to leverage (Table 2).

In Level II, we observe that both `ChatGPT` and `GPT-4` still perform relatively strongly on Derivatives, Corporate Issuers, and Equity Investments, yet still relatively poorly on Financial Reporting. However, the results are now more nuanced: `ChatGPT` struggles on Alternative Investments and Fixed Income compared to `GPT-4`, while `ChatGPT` outperforms `GPT-4` in Portfolio Management and Economics. Interestingly enough, both models now demonstrate low answer accuracy in the Ethics item-sets of Level II. This could originate from the more in-depth, situational, and detailed character of the problems from Level II in comparison to Level I.

CoT prompting yields limited improvements over ZS: Although CoT performs better than ZS in almost all cases and better than FS in Level II for `GPT-4`, we note that the use of CoT did not help LLMs as much as we initially expected (Table 1, Table 3, Table 4). In Level I, CoT prompting hardly benefits `GPT-4` (bringing in just a 1% relative increase) and actually deteriorates the performance of `ChatGPT`. In Level II, CoT prompting yields a decent 7% relative improvement over

ZS prompting for GPT-4, but a disappointing 1% for ChatGPT. Section 5.1 further investigates the reasons explaining such observations. In Level I, we see that CoT negatively affected both LLMs particularly in Quantitative Methods, which could be due to hallucinations in mathematical formula and calculations. In Level II, we notice that CoT benefited both LLMs in the Ethics and Portfolio Management topics, where explicit step-by-step reasoning over long and intricate evidence is usually helpful. In both levels, we also noted that CoT prompting sometimes led to inconsistent performance across questions from the same topic, as manifested by the high standard deviations reported in Table 3 and Table 4.

However, despite the aforementioned observations, it is hard to clearly identify more topics that systematically benefit or suffer from the use of CoT for both models across levels. For instance, in Financial Reporting problems from Level II, GPT-4 saw its accuracy improve by 23% with CoT relative to ZS, while ChatGPT saw its performance decrease by 21% (Table 4).

A few in-context exemplars help more than CoT: Compared to ZS and CoT prompting, FS prompting offers significant performance improvements for ChatGPT on the Level I mock exams (Table 1). 2S prompting yielded the best performance across all categories and overall in Level I for ChatGPT. Across mock exams in Level II, the dominance is not as significant, but FS prompting still manages to achieve the best overall score for both models, with the exception of Level II for GPT-4 (Table 3, Table 4). Interestingly, for Level II, the best FS prompting type was 10S prompting for ChatGPT, which suggests more complex exams benefited from a more holistic FS approach across multiple topics. The overall trend shown in the results is that FS prompting seems to offer better assistance to less complex models (ChatGPT) when being tested on seemingly simpler exams (Level I).

It is likely that FS yields better performance improvement than CoT because it shows actual correct answers to different types of mock questions. It also enables the models to understand how to best use the table evidence or other information contained in a question (if any). The comparatively lower performance improvement brought by FS observed in Level II mock exams may be due

to the more complex nature of the questions and the fact they include case studies; it may be a scenario where simply prompting the models with the correct answers is not sufficient. Level II may thus benefit from a combination of FS and CoT prompting with clear explanations as to how the information in the case study was leveraged to arrive at the correct answer.

5 Detailed Analysis

5.1 Underperformance of CoT on Level I

It was surprising to see that CoT only marginally improved the models' performance on each test, and was actually slightly detrimental to the performance of ChatGPT on the Level I exams. To inspect the nature of the errors made by the models when using CoT prompting, we looked over each instance where no-CoT was correct while CoT was incorrect, and categorized the error as one of: Knowledge, Reasoning, Calculation, or Inconsistency.

Knowledge errors are those where the model lacks critical knowledge required to answer the question. This includes an incorrect understanding of some concept, not knowing the relationship between concepts, or using an incorrect formula to answer a question requiring calculation. Reasoning errors are when the model had all the correct knowledge, but either over-reasoned in its response, or hallucinated some additional requirements or information in the question that was not actually present. Calculation errors are errors pertaining to some incorrect calculation (using a correct formula), or failing to accurately compare or convert results. Errors of inconsistency are when the model's thinking is entirely correct, yet it chooses the wrong answer.

Type of Error	ChatGPT	GPT-4
Knowledge	55.2%	50.0%
Reasoning	8.6%	10.7%
Calculation	17.2%	28.6%
Inconsistency	19.0%	10.7%

Table 5: Error modes of level I questions ChatGPT and GPT-4 got correct without CoT but incorrect using CoT

ChatGPT: By far the most common error mode for ChatGPT is knowledge based, constituting

over half of all errors VS. no-CoT. This implies that, with CoT reasoning, the gaps in the LLMs internal knowledge are magnified. As the model begins to think through its answer, it states its incorrect assumptions, which it proceeds to rationalize in the context of the question thereby skewing the rest of the answer towards a wrong choice. Without using CoT reasoning, the model is able to make an "educated guess" where any incorrect knowledge has less of an opportunity of skewing the guess towards an incorrect answer. With a 1/3 chance of guessing correctly, plus any contextual hints that may lie in the question, for questions where GPT simply lacks the knowledge to reason correctly, guessing is a more accurate strategy.

This same principal similarity explains calculation and reasoning errors, where one or a few off-track token generations then throw off the rest of the answer, resulting in an incorrect conclusion.

The instances where the model is entirely correct but then concludes or just selects the wrong answer are more enigmatic. In about half of these cases, it seems to fail to generate a stop token upon coming to the conclusion, leading it to restate the concluding sentence with another option selected. In the other cases, there appears to be some disconnect between the thought process and the answer selection. As we were using OpenAI's functions API to retrieve structured output, our leading suspicion is that in these cases the ordering outlined in the system prompt was missed or ignored, and the answer was generated first.

GPT-4: There were about half as many instances of CoT making an error not made without CoT for GPT-4, compared to ChatGPT. On these questions, GPT-4 also displays knowledge errors as the most common error mode. However, unlike ChatGPT, almost none of these knowledge errors were using the incorrect formula. This, along with the fact that there were less knowledge errors in total, shows that GPT-4 has more complete internal knowledge of both financial information and especially financial formulas and calculation methods. Rather than knowledge errors, GPT-4's most common error mode on questions requiring calculation are calculation errors. ChatGPT also frequently made these sorts of errors in conjunction with using the wrong formula, which underlines the well-known and more foundational shortcoming of language models' mathematical abilities (Frieder et al., 2023).

GPT-4 also displayed far fewer inconsistency errors than ChatGPT. It appears to have a much stronger ability to connect its reasoning to the answers and to make comparisons. The one error type that GPT-4 makes more frequently than ChatGPT was reasoning errors. It would seem that, along with GPT-4's greater ability to reason, it has a greater chance of "talking itself" into incorrect lines of reasoning.

Type of Error	ChatGPT	GPT-4
Knowledge	70%	80%
Reasoning	20%	20%
Out of Tokens	10%	0%

Table 6: Error modes of level II questions ChatGPT and GPT-4 got correct without CoT but incorrect using CoT

5.2 CoT Benefits on Level II

If CoT amplifies the effect of missing knowledge, and allows LLMs room to miscalculate or "talk themselves" into a wrong answer, one might question why it seemed to help much more on Level II exams. The Level II exam questions require more interpretation of the information, as one needs to figure out what is relevant from the case, and some information may be missing but is expected to be known and needed to answer the question. Using CoT helps the model to reason over the information and filter what is relevant to the question from the case.

5.3 Can LLMs pass the CFA exam?

5.3.1 CFA Level I Passing Score

The CFA Institute refrains from disclosing the minimum passing score (MPS) for its examinations, thereby giving rise to an entire industry centered around speculating on the elusive actual MPS. The MPS is uniquely established for each individual exam, guided by the standards that the CFA Institute established back in 2011.

The CFA Institute employs the 'Angoff Standard Setting Method' to ascertain the pass rates for CFA exams. This involves a group of CFA Charterholders convening to collectively assess the true difficulty level of the questions and the appropriate level of ease that should accompany passing each question.

Exam	ChatGPT			GPT-4		
	ZS	CoT	FS	ZS	CoT	FS
Level I	Pass	Fail	Pass	Pass	Pass	Pass
Level II	Fail	Fail	Fail	Unclear	Pass	Pass

Table 7: ChatGPT and GPT-4 ability to pass Level I and Level II Exams

Although the CFA Institute maintains an air of secrecy surrounding its pass/fail thresholds, certain indicators point towards a potential elevation of the MPS for CFA Level I. Drawing from feedback provided by CFA exam takers on Reddit, the average MPS stood at 65% in December 2019, but surged to 71.1% by February 2021. In June 2019, estimations suggest that certain individuals managed to pass CFA Level I with a mere 60.8%; by February 2021, this had escalated to 66.7%.

Aiming for approximately 70% in as many subjects as possible seems to be a prudent strategy for clearing CFA Level I. Put differently, attaining scores above 70% in all topics is not a necessity for passing. Some contend that achieving as low as 65% or even 63% might suffice. Remarkably, one doesn't even need to exceed 51% in every area to secure a passing grade. The pattern appears to allow for the possibility of scoring below 50% in about three, or perhaps four, subjects. However, this would likely necessitate counterbalancing with scores exceeding 70% in at least three subjects and falling between 51% and 70% in the remaining ones. Nevertheless, maintaining an average score of 70% across subjects considerably enhances the likelihood of a positive outcome upon receiving the results.⁴

5.3.2 CFA Level II Passing Score

The estimations from the Reddit community regarding the MPS for CFA Levels II and III are even more outdated than those for Level I, yet they indicate that the two advanced exams have consistently featured lower passing thresholds. In June 2019, their approximations pegged the MPS for Level III at a mere 57.4%, and for Level II at just 62.8%. The subject level passing scores are ambiguous for the Level II exam, but we can attempt to apply the same logic as the Level I exam but

⁴<https://www.efinancialcareers.com.au/news/finance/whats-the-minimum-score-you-can-get-on-cfa-level-i-and-still-pass>

make an assumption that threshold for each subject is 60% instead of 70%.⁵

5.3.3 Proposed pass criteria

Given the above information our proposed pass criteria is as follows:

- Level I - achieve a score of at least 60% in each topic and an overall score of at least 70%
- Level II - achieve a score of at least 50% in each topic and an overall score of at least 60%

Table 7 shows which model implementations were able to pass the exams. The FS implementations in both settings correspond to the number of shots shown in Table 3 and Table 4. Most of the settings showed a clear pass or fail except for GPT-4 ZS on Level II which was a borderline decision either way. GPT-4 in a ZS setting attains a score of >60% in six of the topics and achieves a score of between 50% and 60% in one of the topics. The topic performance seems high but the overall score of 57.39% falls slightly short of the minimum passing score proposed earlier, it is thus unclear as to whether this LLM setting would pass the CFA Level II exam.

6 Conclusion and Discussion

In this paper, we have conducted a thorough evaluation of ChatGPT and GPT-4 on the CFA level I and level II exams. We observed that GPT-4 performed better than ChatGPT in almost every topic of both levels when using the same prompting method. Based on estimated pass rates and average self-reported scores, we concluded that ChatGPT would likely not be able to pass the CFA level I and level II under all tested settings, while GPT-4 would have a decent chance of passing the CFA Level I and Level II if prompted with FS and/or CoT.

⁵<https://www.efinancialcareers.com.au/news/finance/whats-the-minimum-score-you-can-get-on-cfa-level-i-and-still-pass>

We noted that CoT prompting provided little improvement for ChatGPT on both exams and GPT-4 on the Level I exam. While CoT prompting did help the models reason and understand the question and information better, it also exposed them to making errors due to incorrect/missing domain specific knowledge as well as reasoning and calculation errors. Additionally, we noticed that FS helped LLMs the most in both Levels thanks to the integration of positive instances into the prompt, yielding the best performance in most cases.

With these observations in mind, we propose future systems that could display greater performance by utilizing various tools. The most prevalent error mode of CoT, knowledge errors, could be addressed through retrieval-augmented generation using an external knowledge base containing CFA-specific information. Calculation errors could be avoided by offloading calculations to a function or API such as Wolfram Alpha. The remaining error modes, reasoning and inconsistency, could be reduced by employing a critic model to review and second guess the thinking before submitting the answer, or combining FS and CoT together to give richer examples of expected behavior. We hope this work paves the way for future studies to continue enhancing LLMs for financial reasoning problems through rigorous evaluation.

Acknowledgments

This research was funded in part by the Faculty Research Awards of J.P. Morgan AI Research. The authors are solely responsible for the contents of the paper and the opinions expressed in this publication do not reflect those of the funding agencies.

Disclaimer This paper was prepared for informational purposes by the Artificial Intelligence Research group of JPMorgan Chase & Co and its affiliates (“JP Morgan”), and is not a product of the Research Department of JP Morgan. JP Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solici-

tation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

References

- [Abdaljalil and Bouamor2021] Samir Abdaljalil and Houda Bouamor. 2021. An exploration of automatic text summarization of financial reports. In *Proceedings of the Third Workshop on Financial Technology and Natural Language Processing*, pages 1–7.
- [Araci2019] Dogu Araci. 2019. Finbert: Financial sentiment analysis with pre-trained language models. *arXiv preprint arXiv:1908.10063*.
- [Bang et al.2023] Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity.
- [Barros et al.2023] Thierry S Barros, Carlos Eduardo S Pires, and Dimas Cassimiro Nascimento. 2023. Leveraging bert for extractive text summarization on federal police documents. *Knowledge and Information Systems*, pages 1–31.
- [Brown et al.2020] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- [Chen et al.2022] Zhiyu Chen, Shiyang Li, Charese Smiley, Zhiqiang Ma, Sameena Shah, and William Yang Wang. 2022. Convfinqa: Exploring the chain of numerical reasoning in conversational finance question answering.
- [Chowdhery et al.2022] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways.
- [Frieder et al.2023] Simon Frieder, Luca Pinchetti, Alexis Chevalier, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz, Philipp Christian Petersen, and Julius Berner. 2023. Mathematical capabilities of chatgpt.
- [Katz et al.2023] Daniel Martin Katz, Michael James Bommarito, Shang Gao, and Pablo Arredondo. 2023. Gpt-4 passes the bar exam. *Available at SSRN 4389233*.

- [Khashabi et al.2020] Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Hannaneh Hajishirzi. 2020. Unifiedqa: Crossing format boundaries with a single qa system.
- [Kim et al.2022] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoon Yun, Dongyoon Han, and Seunghyun Park. 2022. Ocr-free document understanding transformer. In *European Conference on Computer Vision*, pages 498–517. Springer.
- [Kung et al.2023] TH Kung, M Cheatham, A Medinilla, C Sillos, L De Leon, C Elepaño, et al. 2023. Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models. *plos digit health* 2 (2): e0000198.
- [Kııcıman et al.2023] Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. 2023. Causal reasoning and large language models: Opening a new frontier for causality.
- [Le et al.2022] Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu Hong Hoi. 2022. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. *Advances in Neural Information Processing Systems*, 35:21314–21328.
- [Lewis et al.2020] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- [Li et al.2022] Jingye Li, Hao Fei, Jiang Liu, Shengqiong Wu, Meishan Zhang, Chong Teng, Donghong Ji, and Fei Li. 2022. Unified named entity recognition as word-word relation classification.
- [Li et al.2023] Xianzhi Li, Xiaodan Zhu, Zhiqiang Ma, Xiaomo Liu, and Sameena Shah. 2023. Are chatgpt and gpt-4 general-purpose solvers for financial text analytics? an examination on several typical tasks. *arXiv preprint arXiv:2305.05862*.
- [Mosbach et al.2023] Marius Mosbach, Tiago Pimentel, Shauli Ravfogel, Dietrich Klakow, and Yanai Elazar. 2023. Few-shot fine-tuning vs. in-context learning: A fair comparison and evaluation. *arXiv preprint arXiv:2305.16938*.
- [Rahimi et al.2023] Saba Rahimi, Tucker Balch, and Manuela Veloso. 2023. Exploring the effectiveness of gpt models in test-taking: A case study of the driver’s license knowledge test. *arXiv preprint arXiv:2308.11827*.
- [Strong et al.2023] Eric Strong, Alicia DiGiammarino, Yingjie Weng, Preetha Basaviah, Poonam Hosamani, Andre Kumar, Andrew Nevins, John Kugler, Jason Hom, and Jonathan Chen. 2023. Performance of chatgpt on free-response, clinical reasoning exams. *medRxiv*, pages 2023–03.
- [Touvron et al.2023] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- [Vaswani et al.2017] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- [Wang et al.2022] Bin Wang, Jiangzhou Ju, Yunlin Mao, Xin-Yu Dai, Shujian Huang, and Jiajun Chen. 2022. A numerical reasoning question answering system with fine-grained retriever and the ensemble of multiple generators for finqa. *arXiv preprint arXiv:2206.08506*.
- [Wang et al.2023] Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2023. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*.
- [Wei et al.2022] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- [Wu et al.2023] Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564*.
- [Yang et al.2023a] Guang Yang, Yu Zhou, Xiang Chen, Xiangyu Zhang, Tingting Han, and Taolue Chen. 2023a. Exploitgen: Template-augmented exploit code generation based on codebert. *Journal of Systems and Software*, 197:111577.
- [Yang et al.2023b] Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. 2023b. Fingpt: Open-source financial large language models. *arXiv preprint arXiv:2306.06031*.
- [Zhang et al.2023] Haopeng Zhang, Xiao Liu, and Jiawei Zhang. 2023. Extractive summarization via chatgpt for faithful summary generation. *arXiv preprint arXiv:2304.04193*.

Appendix

A Topic Distribution in each Level

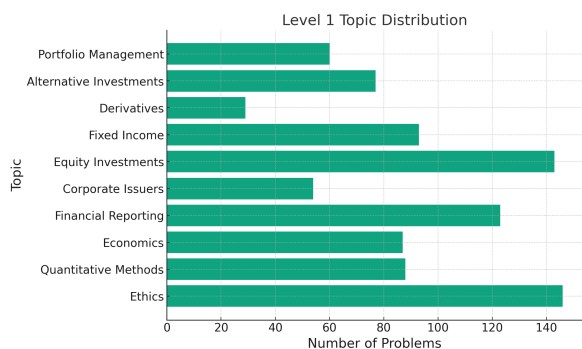


Figure 2: Level I exam topic distribution

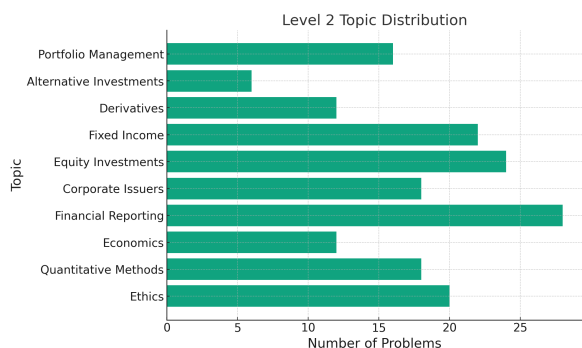


Figure 3: Level II exam topic distribution

B Prompt templates used

B.1 Level I

Listing 1: ZS

SYSTEM: You are a CFA (chartered financial analyst) taking a test to evaluate your knowledge of finance. You will be given a question along with three possible answers (A, B, and C).

Indicate the correct answer (A, B, or C).

USER: Question:
{question}
A. {choice_a}
B. {choice_b}
C. {choice_c}

Listing 2: CoT

SYSTEM: You are a CFA (chartered financial analyst) taking a test to evaluate your knowledge of finance. You will be given a question along with three possible answers (A, B, and C).

Before answering, you should think through the question step-by-step. Explain your reasoning at each step towards answering the question. If calculation is required, do each step of the calculation as a step in your reasoning.

Indicate the correct answer (A, B, or C).

USER: Question:
{question}
A. {choice_a}
B. {choice_b}
C. {choice_c}

Listing 3: FS (2S example)

SYSTEM: You are a CFA (chartered financial analyst) taking a test to evaluate your knowledge of finance. You will be given a question along with three possible answers (A, B, and C).

Indicate the correct answer (A, B, or C).

USER: Question:
{question}
A. {choice_a}
B. {choice_b}
C. {choice_c}

ASSISTANT: {answer}

USER: Question:
{question}
A. {choice_a}
B. {choice_b}
C. {choice_c}

ASSISTANT: {answer}

USER: Question:
{question}
A. {choice_a}
B. {choice_b}
C. {choice_c}

B.2 Level II

For Level II, the case description of each item-set was inserted before each question from the user.