000FAMMA:ABENCHMARKFORFINANCIAL001MULTILINGUALMULTIMODALQUESTION003ANSWERINGOULTIMODALQUESTION

Anonymous authors

006

007

012 013

014

015

016

017

018

019

021

023

025

026

028

029

032

Paper under double-blind review

Abstract

In this paper, we introduce FAMMA, an open-source benchmark for financial multilingual multimodal question answering (QA).¹ Our benchmark aims to evaluate the abilities of multimodal large language models (MLLMs) in answering questions that require advanced financial knowledge and sophisticated reasoning. It includes 1,758 meticulously collected question-answer pairs from university textbooks and exams, spanning 8 major subfields in finance including corporate finance, asset management, and financial engineering. Some of the QA pairs are written in Chinese or French, while a majority of them are in English. These questions are presented in a mixed format combining text and heterogeneous image types, such as charts, tables, and diagrams. We evaluate a range of state-of-the-art MLLMs on our benchmark, and our analysis shows that FAMMA poses a significant challenge for these models. Even advanced systems like GPT-40 and Claude-35-Sonnet achieve only 42% accuracy. Additionally, the open-source Qwen2-VL lags notably behind its proprietary counterparts. Lastly, we explore GPT o1-style reasoning chains to enhance the models' reasoning capabilities, which significantly improve error correction. Our FAMMA benchmark will facilitate future research to develop expert systems in financial QA. The code and data have been anonymously released at https: //github.com/random2024G0/bench-script.

031 1 INTRODUCTION

Benchmarks have played a pivotal role in advancing AI research, particularly in the realm of large
language models (LLMs) (Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023; Jiang et al., 2023;
2024; Meta, 2024). Benchmarks have been helping researchers track the advancement of LLMs in a
variety of capabilities, including general language understanding and knowledge acquisition (Wang
et al., 2019; Hendrycks et al., 2021a; Zhou et al., 2023; Wang et al., 2024b), code generation (Chen
et al., 2021a; Liu et al., 2023; Jimenez et al., 2024), mathematical reasoning (Cobbe et al., 2021;
Hendrycks et al., 2021b), tool use (Yan et al., 2024; Srinivasan et al., 2023; Trivedi et al., 2024), and
legal reasoning (Guha et al., 2023). Meanwhile, we have seen a scarcity of high-quality benchmarks
in financial reasoning, an area where practitioners are eager to benefit from LLMs.

042 We envision that LLMs will have a broad and significant impact in the finance industry, enabling intelligent systems that can assist human experts in various tasks such as risk management and pre-043 dictive analytics. Towards this goal, high-quality benchmarks are needed to track the capabilities of 044 LLMs in understanding financial knowledge and answering complex financial questions. Unfortu-045 nately, existing benchmarks in this domain cannot fully reflect the nature of daily work of financial 046 practitioners: they only have text-based questions; they are all in English; answering their questions 047 only requires knowledge at a rudimentary to intermediate level (Hendrycks et al., 2021a; Chen et al., 048 2021b; Islam et al., 2023; Xie et al., 2024). For a detailed comparison, refer to section 2. However, financial practitioners often have to handle other data modalities, read documents in other languages, and use advanced knowledge. For example, traders often rely on charts to identify trading opportu-051 nities; financial analysts need to analyze documents that include many complicated tables; investors

052

¹This name is made to honor Eugene Fama, a Nobel prize winner that is best known for his work on portfolio theory, asset pricing, and the efficient-market hypothesis. He is regarded as "the father of modern finance".



Figure 1: Sample question in portfolio management, classified as hard difficulty. Typically covered in master's courses. Prerequisites include master's-level knowledge of statistics and optimization theory. See the course description of Fin-504,505 in Princeton and outlines of portfolio management of CFA-Level III.

078

079

sometimes have to read earning reports in languages other than English; quantitative researchers often need to use advanced knowledge such as stochastic calculus to price financial contracts.

In this paper, we present FAMMA, an open-source benchmark for financial multilingual multimodal question answering (QA). Figure 2 displays three QA examples in our benchmark. Compared to existing benchmarks, FAMMA has a significantly better reflection of the real problems that financial practitioners address on a daily basis. This benchmark includes 1,758 meticulously collected question-answer pairs from university textbooks and exams, spanning 8 major subfields including corporate finance, asset management, and financial engineering. Answering these questions requires advanced knowledge such as factor models, option pricing, and asset allocation. A good portion of the QA pairs are written in Chinese or French, although a majority of them are in English. These questions are presented in a mixed format combining text and heterogeneous image types, such as charts, tables, and diagrams.

We evaluate 3 advanced proprietary MLMMs such as GPT-40 (OpenAI, 2023) and 1 top-ranked open-source MLMM—Qwen2-VL (Wang et al., 2024a). Our key findings are summarized below:

- FAMMA presents significant challenges: GPT-40 and Claude-35-Sonnet only achieve 42% accuracy, notably lower than human performance 56%, indicating substantial room for improvement. In addition, there is a pronounced disparity in performance between Qwen2-VL and GPT-40.
- Our analysis of 100 GPT-40 error cases reveals that 42.5% are due to domain knowledge gaps, while 27.5% involve ambiguous responses. This suggests that GPT-40 struggles with financial knowledge and at times generates imprecise answers despite correctly understanding the problem.
- We explore GPT o1-style reasoning chains to enhance the models' reasoning capabilities, significantly outperforming the RAG method in correcting errors on FAMMA, particularly in the categories of ambiguous answer generation and numerical inaccuracy.
- 100 101 102

092

094

096

098

099

2 RELATED WORK

103 104

The application of natural language technologies in finance dates back to the early 2000s, when
 sentiment analysis was used to analyze how media would impact stock market movements (Tetlock,
 2007; Pang et al., 2008). Over recent years, the emergence of LLMs has inspired research in advancing financial industry with LLMs, including pretraining and fine-tuning LLMs with finance-related



(a) Sample question in financial statement analysis, classified as medium difficulty. Typically covered in master's courses. Prerequisites include senior undergraduate's-level or higher knowledge of accounting and corporate finance. See the course description of Fin-502 in Princeton and outlines of financial statement analysis of CFA-Level II.

Context: We have a three-period binomial model shown below. At time zero, we have a stock whose price per share we denote by \mathcal{G}_{n} , a positive quantity known at time zero. At time one, the price per share of this stock will be one of two positive values, where the H and T standing for head and tail, respectively... (details omitted)



Question 1: Assume risk-neural probability for the up and down move are both 0.5, compute the conditional expectation of based on the information at time 1 under the risk neural measure $E[S_2](H), E[S_2](T)$

Question 2: Under the actual probability with the up and down move probability 2/3 and 1/3, consider the maximum-to-date process below $M_n = \max_{max} S_n$, compute $E_2[M3](TH) E_2[M3](TT)$ and determine whether M_n is Markov?

(b) Sample question in derivatives, classified as hard difficulty. Typically covered in master's courses. Prerequisites include master's level or higher knowledge of probability theory and stochastic calculus. See the course description of Fin-501,503 in Princeton and outlines of derivatives of CFA-Level III.

- Figure 2: Sampled FAMMA examples from the other two subfields. The questions and images need expertlevel knowledge to understand and reason. Samples in this figure are text truncated due to space.
- 135 136

text (Wu et al., 2023; Yang et al., 2023), improving sentiment analysis with LLMs (Konstantinidis
et al., 2024; Inserte et al., 2024; Cao et al., 2024), and building chatbots that specialize in finance
knowledge (Chase, 2022; Stratosphere-Technology, 2023; Xue et al., 2023; 2024).

140 Several existing benchmarks can be used to evaluate these modern models and systems, includ-141 ing FiQA (Maia et al., 2018), FinQA (Chen et al., 2021b), ConvFinQA (Chen et al., 2022), Fi-142 nanceBench (Islam et al., 2023), and FinBen (Xie et al., 2024). However, these benchmarks cannot 143 reflect the nature of real problems that financial practitioners have to deal with on a daily basis. In particular, their data only has text but not data of other modalities; their data is only in English; their 144 questions only test knowledge at a rudimentary to intermediate level. The finance-related questions 145 in MMMU (Yue et al., 2024) involve data of other modalities like tables and charts. However, this 146 general-purpose benchmark covers multiple disciplines (e.g., art, business, science, humanities, etc), 147 and thus has a very limited coverage on finance-related questions. Our FAMMA benchmark makes 148 a unique and focused contribution to the community on top of existing benchmarks: it has a much 149 broader coverage on subfields of finance; its data is in multiple languages and of multiple modalities; 150 its questions test advanced knowledge. 151

152 153

154

3 THE FAMMA BENCHMARK

FAMMA provides comprehensive coverage across eight key subfields: alternative investments, corporate finance, derivatives, economics, equity, financial statement analysis, fixed income, and portfolio management. These topics closely align with those taught in elite academic programs, such as Princeton's Master in Finance, as well as professional certifications like the CFA program. The dataset consists of both multiple-choice (55.5%) and open questions (45.5%). Additionally, 70.4% of the questions feature single-image scenarios, while 29.6% involve multi-image scenarios. Notably, 99.5% of the questions are accompanied by explanations. The questions are distributed across three difficulty levels and three most widely used languages in the finance industry (eFinancialCa-

162 reers, 2022): English (78.8%), Chinese (14.4%), and French (6.8%). FAMMA is divided into a 163 validation set and a test set. The validation set, useful for hyperparameter selection, contains 120 164 questions, while the test set comprises 1638 questions The overall subject coverage and statistics are 165 shown in Table 1 while distribution of questions by languages and subfields are shown in Figure 3 166 and Figure 4, respectively. More detailed descriptive statistics can be found in Table 7 and Table 8 in Appendix A. 167

TOTAL QUESTIONS * MULTIPLE-CHOICE * OPEN	1758 976 (55.5%) 782 (44.5%)
* MULTIPLE-CHOICE * OPEN	976 (55.5%) 782 (44.5%)
* Open	782 (44.5%)
	200 / 100 / - 10
DIFFICULTIES	608 / 438 / 712
(Easy:Medium:Hard)	34.6%:24.9%:40.5%
	# TOKENS AVG
By INPUT AND OUTPUT	
* QUESTIONS	233.43
* EXPLANATION	73.95
BY SPLITS	
* VALIDATION	224.29
* TEST	234.12
BY LANGUAGES	
* English	257.11
* CHINESE	94.06
* FRENCH	254.28
By DIFFICULTY LEVELS	
* EASY	136.62
* Medium	109.68
* HARD	375.25
BY SUBFIELDS	
ALTERNATIVE INVESTMENTS	473.33
CORPORATE FINANCE	81.74
* DERIVATIVES	277.47
* ECONOMICS	567.64
* EQUITY	243.15
* FINANCIAL STATEMENT ANALYSIS	49.67
* FIXED INCOME	198.21
PORTFOLIO MANAGEMENT	276.19

Table 1: Key statistics of FAMMA.



English Chinese French

Figure 3: Distribution of questions in FAMMA across languages.



Figure 4: Distribution of questions in FAMMA across subfields.

199 3.1 DATASET CONSTRUCTION 200

201 Question collection. We assembled a team of seven volunteer STEM researchers to create a com-202 prehensive set of multimodal questions. Five are co-authors of this paper, while the other two are 203 graduates from a Chinese university. Two annotators hold finance degrees, and the others have com-204 pleted relevant coursework. These annotators draw upon open-source textbooks, exams, and other 205 study materials (see Table 6 in Appendix A for details), and apply their expertise to rewrite or create 206 new questions when needed. The new questions are either entirely original, not present in the data 207 sources, or enhanced versions of existing questions.

208 The annotators are tasked with selecting questions that require advanced, master-level, or profes-209 sional knowledge to answer. This selection process is guided by aligning the questions with a min-210 imum of CFA Level 1 difficulty (CFA Institute, 2024b), ensuring they meet industry standards of 211 complexity. Additionally, selected questions must incorporate multimodal information, such as ta-212 bles, images, or other visual data, to enrich the input and challenge the model's ability to process diverse formats. By following these criteria, we have curated a diverse set of approximately 2,000 213 questions, drawn from a wide range of authoritative sources. 214

215

168

197

Data quality control. We follow a two-stage data cleaning process to ensure the data quality.

In the first stage, we conduct a thorough review to correct formatting errors, fix typos, remove duplicate questions, and verify the accuracy of explanations. Each question is cross-verified by 2-3 annotators to ensure consistency and accuracy. Formatting errors and typos arise due to variations in the original sources, such as UTF encoding issues in Chinese and French texts, and the explanations are either provided by the source materials or written by annotators.

In the second stage, we classify each question into one of three difficulty levels—easy, medium, or hard—and label it with the appropriate subfield.

The difficulty levels are aligned with the concept-specific standards of the CFA curriculum (CFA 224 Institute, 2024a). In addition, questions that require processing more complex information, such as 225 multiple tables and images, are considered more difficult. In cases where the difficulty is ambiguous, the annotators use their judgment. Additionally, questions deemed overly simplistic—such 226 as those based purely on memorization or with answers that are obvious from the context-are 227 removed to maintain the desired level of challenge and to ensure they test knowledge and reason-228 ing. The subfield annotation is determined by the explicit topics provided in the data source. If the 229 subfield is not clearly specified, the annotators use their discretion to assign the most appropriate 230 category based on the content of the question. 231

The JSON formats of the multiple-choice and open questions are illustrated in Listing 5 and Listing 6 in Appendix A, respectively.

To conclude, the FAMMA dataset offers a diverse range of questions, enabling the evaluation of models across various scenarios and allowing for fine-grained analysis of their performance.

4 EXPERIMENTS

237

238 239

240 241

242

243

244

245 246

247 248

249

260

4.1 EXPERIMENTAL SETUP

Benchmarked MLMMs. We evaluate three cutting-edge closed-source models that are ranked among the top 10 on the Multimodal Arena Leaderboard (Lmsys Org, 2024):

- GPT-40 (OpenAI, 2024a): The latest iteration in the GPT series, GPT-40 features enhanced capabilities in language and vision understanding, as well as improved generation performance.
- Claude-3.5-Sonnet (Anthropic, 2024): Developed by Anthropic, Claude Sonnet introduces architectural innovations that improve multimodal dialogue and reduce harmful outputs.
- Claude-3-Sonnet (Anthropic, 2024): An earlier version of the Claude Sonnet model.

Additionally, we assess a leading open-source model: Qwen2-VL (Yang et al., 2024) that achieves
state-of-the-art performance on visual understanding benchmarks, including MathVista (Lu et al., 2024) and DocVQA (Mathew et al., 2021).

Generation process. MLMMs are instructed to understand the format and the structure of the questions, and return the response, under a zero-shot setting on our benchmark. The instruction prompts are designed to be straightforward and consistent across all models. Please refer to Listing 7 and Listing 8 in Appendix B for the prompts used to to guide responses to multiple-choice and open questions, respectively. During the final stage of generation for multiple-choice questions, we utilize both regex and GPT-40 to extract the corresponding lettered option from the response. Any discrepancies between the two methods will be manually reviewed and validated by annotators.

LM-powered evaluation. During the evaluation process, we use GPT-40 as an LM evaluator to assess the accuracy of responses generated by LLMs for each question. The reported score represents the accuracy of these responses. Each response is categorized as either correct or incorrect, and the reported score reflects the average accuracy across the entire set of questions.

The LM evaluator is instructed to understand the format and structure of the questions, as well as to consider the key points in the ground-truth answers for evaluating the responses. Please refer to Listing 9 in Appendix B for the instructions provided for evaluating the answers. Note that for open-ended questions, where both gold and generated answers are provided, there is a single correct answer, making the 1-0 correctness straightforward to determine. We set the temperature of the LM evaluator as 0 to keep the evaluation results deterministic. Human performance. We invite two volunteers to participate in the test to establish a human benchmark. Both are experienced finance professionals: one holds a Master's degree in Finance from a U.S. institution, while the other graduated from a Grande École in France, specializing in mathematics and finance. The first volunteer, proficient in both English and Chinese, is tasked with completing half of the English test and the entire Chinese test, while the second volunteer takes the remaining portion of the English test and the full French test. They are allowed to consult textbooks, e.g. Hull (2017); Bodie et al. (2014), but are prohibited from searching the web for answers.

It worth noting that the human score is roughly the same as those estimated from CFA passing scores.
Based on the report ², the passing score is approximately 68% for all the three levels. During the annotation process, the difficulty levels of FAMMA's questions—easy, medium, and hard—closely correspond to those of CFA Levels I, II, and III. Based on this data and assumptions about the performance of unqualified candidates from previous levels, we estimate the accuracy rate for easy, medium, and hard questions to be equal to that of Level I, II, III—68%, 62.24%, 57.26%, respectively, which resulting a overall score of 62.1%. See Appendix B.1 for details of the estimation.

284 285

286

287

288

289

4.2 RESULTS AND ANALYSIS

Main results. We repeat the generation and evaluation process three times, and report the average result along with the standard error across all experiments. See the overall scores, breakdown by difficulty levels and languages in Table 2. We summarize the key findings as follows.

- FAMMA presents a comprehensive challenge. Human performance sets the highest benchmark with an overall score of 56.96, leading across all difficulty levels. Among the models, GPT-40 ranks first with a score of 42.11, followed closely by Claude-35-Sonnet at 41.87. Both models fall approximately 15 points short of human performance and, based on our estimates, about 20 points below CFA professional levels. This substantial gap underscores the significant challenges FAMMA poses for MLLMs.
- The open-source Qwen2-VL significantly lags behind more advanced closed-source MLLMs. Ac-296 cording to its technical report (Wang et al., 2024a), Qwen2-VL has not been explicitly optimized 297 for financial corpora, whereas the Claude family models prioritize finance as a key domain for 298 evaluation and improvement (Claude, 2024). Interestingly, Qwen2-VL performs better on hard 299 questions than on medium ones. A possible explanation is that hard questions often require higher 300 computational complexity and advanced mathematical reasoning, areas where Qwen2-VL excels. 301 In fact, its technical report highlights superior performance on MathVista (Lu et al., 2024), out-302 performing other MLLMs, including GPT-40 and the Claude models. 303
- To conclude, the main results highlight the progress in MLMM QA in finance but also underscores the challenge of surpassing human-level performance.
- Analysis I: model performance across different subfields. As shown in Figure 5, GPT-40 demonstrates the largest margin in economics, a social science discipline that studies the behaviour and interactions of economic agents. The result indicates its rich knowledge in social domains in addition to mathematics reasoning. GPT-40 also excels at financial statement analysis, whose context usually contains long tables (see Figure 2a), indicating its superior ability in table understanding (Kim et al., 2024). This well-rounded performance suggests that GPT-40 possesses a broad understanding of diverse financial concepts, excelling in knowledge-based and applied assessments.
- Claude-35-Sonnet leads in corporate finance, alternative investments, derivatives, and fixed income, though with small margins over GPT-40. Both Qwen2 and Claude-3-Sonnet fall significantly short in most areas. The notable improvement in finance-related QA from Claude-3 to Claude-35-Sonnet is consistent with public findings, where the win rate on finance tasks improved by 27%, as reported in the technical report (Claude, 2024).
- 319

306

Analysis II: model performance across different languages. Seen from Table 2, a consistent observation is that all models perform best in English, with GPT-40 and Claude-35-Sonnet comparably surpass the other competitors with a score around 44. For both GPT-40 and Claude-35-Sonnet, Chinese performance falls noticeably behind English, especially in harder categories, suggesting that

²https://300hours.com/cfa-passing-score/

324		~		-		
325	MODELS	SIZE	OVERALL	EASY	MEDIUM	HARD
326	Human	N/A	56.16	61.35	57.09	52.11
327	PERFORMANCE					
328	GPT-40	N/A	42.85	47.25	39.85	40.98
329			(0.45)	(4.72)	(1.71)	(4.59)
330 331	* English		44.90	48.00	43.83	42.71
332			(0.10)	(4.67)	(0.19)	(4.52)
333	* CHINESE		37.70	39.45	37.35	30.65
334			(2.80)	(3.95)	(4.04)	(2.04)
335	* French		32.50	64.75	31.55	22.90
336			(1.65)	(4.17)	(2.63)	(0.73)
337	CLAUDE 25 Sc		12.80	17 27	41.10	20.80
338	CLAUDE-55-50	INNET IN/A	42.80	47.57	41.10	(2.80)
339	* ENGLISH		(0.49)	(4.12)	(2.72)	(5.89)
340	- ENGLISH		44.20	47.51	40.15	40.89
342	* CHINESE		(0.33)	(4.47)	(0.39)	(3.52)
343	CHINESE		37.50	47.57	37.50	(2.15)
344	* Encueu		(2.92)	(3.55)	(4.24)	(2.15)
345	FRENCH		37.50	(2, 27)	(2.02)	40.00
346			(1.15)	(3.27)	(2.93)	(1.69)
347	QWEN2-VL	70B	34.50	38.39	27.40	35.67
348			(0.33)	(2.52)	(3.43)	(4.19)
349	* English		36.20	38.73	28.21	37.08
350			(0.51)	(4.17)	(1.19)	(3.07)
351	* CHINESE		29.60	34.21	31.25	18.37
352			(1.32)	(2.15)	(4.02)	(1.98)
353 354	* French		25.79	50.00	18.42	34.38
355			(1.85)	(3.25)	(2.68)	(2.19)
356	CLAUDE-3-SON	NNET N/A	31.55	31.91	29.00	32.58
357			(0.34)	(4.42)	(2.22)	(4.08)
358	* English		32.70	33.27	31.62	32.65
359			(0.28)	(4.17)	(1.39)	(3.82)
360	* CHINESE		25.32	19.74	28.91	24.49
361			(3.22)	(4.15)	(3.98)	(2.92)
362	* French		30.53	50.00	21.05	43.75
363			(1.28)	(3.77)	(3.13)	(2.82)
364						

366

367

Table 2: The score of various models on the FAMMA test set, with standard errors indicated in parentheses. The anonymous live-updating leaderboard is available at: https://random2024go.github.io/ indexPage/

368 369

370

there are significant gaps in how well these models handle complex tasks in Chinese. The com paratively low scores for Chinese might reflect the challenges related to tokenization or potentially
 smaller and less diverse training corpora in Chinese relative to English. In addition, Qwen2-VL,
 perform poorly on Chinese test, indicating it falls short of training on Chinese financial corpus

Claude-35-Sonnet outperforms its competitors in French, likely due to the Claude family's focus on optimizing non-English languages, such as Spanish and French (Claude, 2023; 2024). Interestingly, it achieves a higher score on hard questions compared to medium ones. It's important to highlight that the number of hard questions in French is limited to just 32 (see Table 8 in Appendix A), primar-



ror types are less frequent but still notable, showing some difficulty in interpreting context or visual

-102		
433	:	:
101		•
434	Small currency trades and small	The expected holding-period return (
435	exchange-traded derivatives trades	HPR) is calculated by multiplying
436	are typically implemented using the	each possible return by its
437	direct market access (DMA) approach,	probability and summing the results.
400	and the high-touch agency approach is	The calculation is as follows: (0.30
430	typically used to execute large,	* 18%) + (0.50 * 12%) + (0.20 * -5%)
439	non-urgent trades in fixed-income and	= 5.4% + 6% - 1% = 10.4%. Therefore,
440	exchange-traded derivatives markets.	the expected HPR for KMP stock is
441	Listing 1. A second second second of CDT 4. due	10.88%.

Listing 1: A sample case response of GPT-40 due to domain knowledge gaps: the high-touch agency approach is mainly used for illiquid orders but exchange-traded derivatives tend to be very liquid.

Listing 2: A sample case response of GPT-40 due to ambiguous answer generation: it correctly performs the computation but reaches to a wrong final result.

data. These findings underscore the need for improvements in finance domain-specific training, and answer generation to enhance GPT-4o's performance.

	TOP 1 ERR	TOP 2 ERR		TOP 1 ERR	TOP 2 ERR
English	DKG (32.65%)	AAG (25.45%)	EASY	AAG (34.33%)	DKG (32.20%)
IESE	DKG (44.16%)	AAG (34.08%)	Medium	AAG (36.49%)	DKG (30.12%)
NCH	DKG (46.08%)	AAG (39.04%)	HARD	DKG (57.08%)	AAG (26.12%)

Table 3: Error types breakdown by language in the sampled set.

Table 4: Error types breakdown by difficulty in the sampled set.

Analysis V: errors across languages and difficulties. In all languages, DKG consistently ac-459 count for the highest proportion of errors, with French leading at 46.08%, indicating GPT-40 strug-460 gles with understanding finance domain-specific knowledge, particularly in non-English contexts. Interestingly, AAG error emerges as the second most common error type across all three languages, 462 suggesting that despite differences in language complexity, GPT-40 often provides unclear or in-463 complete answers regardless of the language. 464

When analyzing error types across difficulty levels, AAG dominates in both easy (34.33%) and 465 medium (36.49%) categories, while DKG takes the lead in hard questions (57.08%). This shift 466 suggests that for easier questions, the model tends to generate ambiguous answers, likely due to 467 overgeneralization or incomplete interpretations. However, as the complexity increases, the GPT-468 40's lack of domain knowledge becomes more evident. 469

Analysis VI: can RAG or o1-reasoning chain help? We explore two independent methods for improving GPT-4o's performance on FAMMA:

• Retrieval augmented generation (RAG): we augment GPT-40 with external financial knowledge base by incorporating content from textbooks "CFA Level III SchweserNotes, Books 1-5, 2023", which comprehensively cover most of the topics included in FAMMA.

• Dynamic Chain-of-Thought (COT) prompting: we implement o1-style reasoning chains (OpenAI, 2024b), where at each step, GPT-40 can either proceed to the next reasoning step (by trying multiple methods, exploring alternative answers, or questioning previous solutions) or provide a final answer. The process begins with a system prompt that includes instructions to guide a step-by-step reasoning approach. Once the problem is introduced as a user message, an assistant message is pre-loaded to establish a standardized starting point for the GPT-4o's response generation.

481 482

120

442

443

444

445 446 447

448

456

457 458

461

470

471

472 473

474

475

476

477

478

479

480

See Appendix B.3 for the implementation of RAG and dynamic COT (prompts in Listing 16-483 Listing 18). These two strategies are applied on the same error set again. Table 5 shows dynamic 484 COT significantly outperforms RAG, correcting 44 errors compared to RAG's 20. It is particularly 485 effective in the AAG category and also excels in NI and ICU corrections. Listing 4 presents an example of the reasoning process, involving identification, calculation, and self-checking, to ultimately arrive at an answer. Listing 3 illustrates the final answer generation, clearly reaching the correct result instead of producing an ambiguous response. This highlights dynamic COT's ability to resolve unclear answers and enhance numerical reasoning. As expected, RAG performs better at addressing DKG errors, though both methods show limited effectiveness in correcting VDM.

	RAG	DYNAMIC COT
# CORRECTED	20	44
- DKG	10	8
- AAG	5	20
- NI	8	13
- VDM	0	0
- ICU	1	3

Table 5: Improvement by applying RAG and dynamic COT in the sampled set.

```
The table provided gives us the
504
         probabilities and holding-period
505
         returns (HPR) for KMP stock under
506
         different economic scenarios: Boom,
         Normal Growth, and Recession.
507
         The expected holding-period return (
         HPR) can be calculated using the
509
         formula for expected value:
510
511
         Boom: Probability = 0.30, HPR = 18\%
512
         n2. Normal Growth: Probability =
513
         0.50, HPR = 12%\n3. Recession:
514
         Probability = 0.20, HPR = -5\%.
515
         To find the expected HPR:\nE(HPR) =
516
         (0.30 * 18\%) + (0.50 * 12\%) + (0.20)
517
         * -5\%) = (0.30 * 18/100) + (0.50 *
518
         12/100 + (0.20 * -5/100)\n= 0.054 +
519
         0.06 - 0.01 \ n= 0.104 \ or \ 10.40\%
520
```

Listing 3: A corrected AAG case using dynamic COT, as shown in Listing 2: no additional incorrect results were generated.

ł "title": "Understanding the Given Probability Distribution" "content": "The table provided gives us the probabilities and holding-period returns (HPR) for KMP stock under different economic scenarios: Boom, Normal Growth, and Recession.", "next_action": "continue" ł "title": "Calculating the Expected Holding-Period Return" "content": "The expected holding -period return (HPR) can be calculated by ... (omitted)" "next_action": "continue" }, "title": "Performing the Expected Value Calculation", "content": "To find the expected HPR:...(omitted)", "next_action": "continue" "title": "Identifying the Correct Answer" "content": "Based on..., the expected holding-period return for KMP stock is 10.40%.", "next_action": "final_answer" ł

Listing 4: Reasoning process using dynamaic COT.

5 CONCLUSION

In this paper, we introduced a comprehensive benchmark for multilingual multimodal QA within the financial domain, addressing the growing need for robust systems that can interpret and respond to queries in various languages while leveraging diverse data modalities.

532

533

521

522

523

524 525 526

527 528

529

530

531

491 492

493

494

495

496

497

498

499

500

501

- 534 535
- 536
- 527
- 538
- 539

540 REFERENCES

554

559

567

572

579

586

- 542 Anthropic. Claude, 2024.
- Bodie, Z., Kane, A., and Marcus, A. J. *Investments*. McGraw-Hill Education, 10th edition, 2014.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam,
 P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R.,
 Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray,
 S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D.
 Language models are few-shot learners. In *Advances in Neural Information Processing Systems*(*NeurIPS*), 2020.
- Cao, Y., Chen, Z., Pei, Q., Dimino, F., Ausiello, L., Kumar, P., Subbalakshmi, K. P., and Ndiaye, P. M. Risklabs: Predicting financial risk using large language model based on multi-sources data. *arXiv preprint arXiv:2401.07452*, 2024.
- 555 CFA Institute. Cfa curriculum exam topics, 2024a.
- 556557 CFA Institute. Cfa level i cfa exam structure, 2024b.
- ⁵⁵⁸ Chase, H. LangChain, 2022.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. D. O., Kaplan, J., Edwards, H., Burda, Y.,
 Joseph, N., Brockman, G., et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021a.
- 563 Chen, Z., Chen, W., Smiley, C., Shah, S., Borova, I., Langdon, D., Moussa, R., Beane, M., Huang,
 564 T.-H., Routledge, B., and Wang, W. Y. Finqa: A dataset of numerical reasoning over financial
 565 data. Proceedings of the Conference on Empirical Methods in Natural Language Processing
 566 (EMNLP), 2021b.
- Chen, Z., Li, S., Smiley, C., Ma, Z., Shah, S., and Wang, W. Y. Convfinqa: Exploring the chain of
 numerical reasoning in conversational finance question answering. *Proceedings of the Conference* on Empirical Methods in Natural Language Processing (EMNLP), 2022.
- ⁵⁷¹ Claude. The claude 3 model family: Opus, sonnet, haiku. *www.anthropic.com*, 2023.
- 573 Claude. Claude 3.5 sonnet model card addendum. *www.anthropic.com*, 2024.
- 574
 575
 576
 576
 576
 577
 576
 577
 576
 577
 578
 579
 579
 570
 570
 570
 571
 572
 574
 574
 574
 574
 575
 576
 577
 576
 577
 577
 578
 578
 578
 579
 579
 579
 570
 570
 570
 571
 572
 574
 574
 574
 575
 576
 577
 577
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
- eFinancialCareers. The top french quant finance programs, 2022. Accessed: 2024-09-23.
- Guha, N., Nyarko, J., Ho, D. E., Ré, C., Chilton, A., Narayana, A., Chohlas-Wood, A., Peters, A., Waldon, B., Rockmore, D. N., Zambrano, D., Talisman, D., Hoque, E., Surani, F., Fagan, F., Sarfaty, G., Dickinson, G. M., Porat, H., Hegland, J., Wu, J., Nudell, J., Niklaus, J., Nay, J., Choi, J. H., Tobia, K., Hagan, M., Ma, M., Livermore, M., Rasumov-Rahe, N., Holzenberger, N., Kolt, N., Henderson, P., Rehaag, S., Goel, S., Gao, S., Williams, S., Gandhi, S., Zur, T., Iyer, V., and Li, Z. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models. *arXiv preprint arXiv:2308.11462*, 2023.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring
 massive multitask language understanding. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021a.
- Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J.
 Measuring mathematical problem solving with the math dataset. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021b.
 - Hull, J. Options, Futures, and Other Derivatives. Pearson, 10th edition, 2017.

- 594 Inserte, P. R., Nakhlé, M., Qader, R., Caillaut, G., and Liu, J. Large language model adaptation for 595 financial sentiment analysis. arXiv preprint arXiv:2401.14777, 2024. 596 Islam, P., Kannappan, A., Kiela, D., Qian, R., Scherrer, N., and Vidgen, B. Financebench: A new 597 benchmark for financial question answering. arXiv preprint arXiv:2311.11944, 2023. 598 Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. d. l., Bressand, 600 F., Lengyel, G., Lample, G., Saulnier, L., et al. Mistral 7b. arXiv preprint arXiv:2310.06825, 601 2023. 602 Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D. S., Casas, 603 D. d. l., Hanna, E. B., Bressand, F., et al. Mixtral of experts. arXiv preprint arXiv:2401.04088, 604 2024. 605 606 Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., and Narasimhan, K. Swe-bench: 607 Can language models resolve real-world github issues? In Proceedings of the International Conference on Learning Representations (ICLR), 2024. 608 609 Kim, Y., Yim, M., and Song, K. Y. Tablevqa-bench: A visual question answering benchmark on 610 multiple table domains, 2024. 611 612 Konstantinidis, T., Iacovides, G., Xu, M., Constantinides, T. G., and Mandic, D. Finllama: Financial 613 sentiment classification for algorithmic trading applications. arXiv preprint arXiv:2403.12285, 2024. 614 615 Liu, J., Xia, C. S., Wang, Y., and Zhang, L. Is your code generated by chatGPT really correct? rigor-616 ous evaluation of large language models for code generation. In Advances in Neural Information 617 Processing Systems (NeurIPS), 2023. 618 Lmsys Org. Multimodal chatbot arena, 2024. 619 620 Lu, P., Bansal, H., Xia, T., Liu, J., Li, C., Hajishirzi, H., Cheng, H., Chang, K.-W., Galley, M., and 621 Gao, J. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. 622 In Proceedings of the International Conference on Learning Representations (ICLR), 2024. 623 Maia, M., Handschuh, S., Freitas, A., Davis, B., McDermott, R., Zarrouk, M., and Balahur, A. 624 Www'18 open challenge: Financial opinion mining and question answering. In Proceedings of 625 the International World Wide Web Conference (WWW), 2018. 626 627 Mathew, M., Karatzas, D., and Jawahar, C. V. Docvqa: A dataset for vqa on document images. In 628 2021 IEEE Winter Conference on Applications of Computer Vision (WACV), 2021. 629 Meta, A. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024. 630 631 OpenAI. Gpt-4v(ision) system card, 2023. 632 OpenAI. Gpt-40 system card, 2024a. 633 634 OpenAI. Gpt-4o1: A large language model for advanced reasoning and understanding. 2024b. 635 Pang, B., Lee, L., et al. Opinion mining and sentiment analysis. Foundations and Trends® in 636 information retrieval, 2008. 637 638 Srinivasan, V. K., Dong, Z., Zhu, B., Yu, B., Mosk-Aoyama, D., Keutzer, K., Jiao, J., and Zhang, J. 639 Nexusraven: a commercially-permissive language model for function calling. In NeurIPS 2023 640 Foundation Models for Decision Making Workshop, 2023. 641 Stratosphere-Technology. FinChat, 2023. 642 643 Tetlock, P. C. Giving content to investor sentiment: The role of media in the stock market. The 644 Journal of finance, 2007. 645 Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., 646
- 647 Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv* preprint arXiv:2307.09288, 2023.

- Trivedi, H., Khot, T., Hartmann, M., Manku, R., Dong, V., Li, E., Gupta, S., Sabharwal, A., and Balasubramanian, N. Appworld: A controllable world of apps and people for benchmarking interactive coding agents. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, 2024.
- Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2019.
- Wang, P., Bai, S., Tan, S., Wang, S., Fan, Z., Bai, J., Chen, K., Liu, X., Wang, J., Ge, W., Fan,
 Y., Dang, K., Du, M., Ren, X., Men, R., Liu, D., Zhou, C., Zhou, J., and Lin, J. Qwen2-vl:
 Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024a.
 - Wang, Y., Ma, X., Zhang, G., Ni, Y., Chandra, A., Guo, S., Ren, W., Arulraj, A., He, X., Jiang, Z., et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. arXiv preprint arXiv:2406.01574, 2024b.
- Wu, S., Irsoy, O., Lu, S., Dabravolski, V., Dredze, M., Gehrmann, S., Kambadur, P., Rosenberg, D., and Mann, G. Bloomberggpt: A large language model for finance. arXiv preprint arXiv:2303.17564, 2023.
- Kie, Q., Han, W., Chen, Z., Xiang, R., Zhang, X., He, Y., Xiao, M., Li, D., Dai, Y., Feng, D., Xu, Y., Kang, H., Kuang, Z., Yuan, C., Yang, K., Luo, Z., Zhang, T., Liu, Z., Xiong, G., Deng, Z., Jiang, Y., Yao, Z., Li, H., Yu, Y., Hu, G., Huang, J., Liu, X.-Y., Lopez-Lira, A., Wang, B., Lai, Y., Wang, H., Peng, M., Ananiadou, S., and Huang, J. The finben: An holistic financial benchmark for large language models. *arXiv preprint arXiv:2402.12659*, 2024.
- Kue, S., Zhou, F., Xu, Y., Jin, M., Wen, Q., Hao, H., Dai, Q., Jiang, C., Zhao, H., Xie, S., He, J.,
 Zhang, J., and Mei, H. Weaverbird: Empowering financial decision-making with large language
 model, knowledge base, and search engine. *arXiv preprint arXiv:2308.05361*, 2023.
- Kue, S., Qi, D., Jiang, C., Shi, W., Cheng, F., Chen, K., Yang, H., Zhang, Z., He, J., Zhang, H., Wei,
 G., Zhao, W., Zhou, F., Yi, H., Liu, S., Yang, H., and Chen, F. Demonstration of db-gpt: Next
 generation data interaction system empowered by large language models. In *Proceedings of the VLDB Endowment*, 2024.
- Yan, F., Mao, H., Ji, C. C.-J., Zhang, T., Patil, S. G., Stoica, I., and Gonzalez, J. E. Berkeley function
 calling leaderboard, 2024.
- Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C., Li, C., Li, C., Liu, D., Huang, F., Dong, G., Wei, H., Lin, H., Tang, J., Wang, J., Yang, J., Tu, J., Zhang, J., Ma, J., Yang, J., Xu, J., Zhou, J., Bai, J., He, J., Lin, J., Dang, K., Lu, K., Chen, K., Yang, K., Li, M., Xue, M., Ni, N., Zhang, P., Wang, P., Peng, R., Men, R., Gao, R., Lin, R., Wang, S., Bai, S., Tan, S., Zhu, T., Li, T., Liu, T., Ge, W., Deng, X., Zhou, X., Ren, X., Zhang, X., Wei, X., Ren, X., Liu, X., Fan, Y., Yao, Y., Zhang, Y., Wan, Y., Chu, Y., Liu, Y., Cui, Z., Zhang, Z., Guo, Z., and Fan, Z. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- Yang, H., Liu, X.-Y., and Wang, C. D. Fingpt: Open-source financial large language models. arXiv preprint arXiv:2306.06031, 2023.
- Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun,
 Y., Wei, C., Yu, B., Yuan, R., Sun, R., Yin, M., Zheng, B., Yang, Z., Liu, Y., Huang, W., Sun,
 H., Su, Y., and Chen, W. Mmmu: A massive multi-discipline multimodal understanding and
 reasoning benchmark for expert agi. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
 - Zhou, J., Lu, T., Mishra, S., Brahma, S., Basu, S., Luan, Y., Zhou, D., and Hou, L. Instructionfollowing evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.

697

698

682

660

661

662

663

702 703	MATERIALS	LANGUAGE	SOURCE
704 705	QUIZZES ON FINANCE-RELATED COURSES	English	MIT OPEN- Course
706	TEXTBOOK: NUMERICAL PROBABILITY	English	ONLINE PDF
707	TEXTBOOK: PAUL WILMOTT ON QUANTITATIVE FINANCE	English	GITHUB
708 709	QUIZZES ON FINANCE-RELATED COURSES	CHINESE	Renming University
710 711 712	QUIZZES ON FINANCE-RELATED COURSES QUIZZES ON FINANCE-RELATED COURSES	French French	Scribd Academia

Table 6: Selected sources as references for generating questions.

Appendices

A DATASET DETAILS

Data source. The question-response pairs are primarily collected from free online resources, quizzes, textbooks, and other study materials. See Table 6 for more details.

Data format. Following data validation, we provide the following information for each question:

- Question ID: a unique identifier for the question across the whole dataset.
- Context: relevant background information related to the question.
- Question: the specific query being asked.
- Images: directories of images referenced in the context or question.
- Options: a list of possible answers, applicable only to multiple-choice questions.
- Question type: categorized as either multiple-choice or open-ended.
- Main question ID: a unique identifier for the question within its context; questions with the same context share the same ID.
- Sub question ID: a unique identifier for the question within its corresponding main question.
- Answer: a concise and accurate response.
 - Explanation: a detailed justification for the answer.
 - Images for explanation: directories of images supporting the explanation.
 - Subfield: the specific area of expertise to which the question belongs, categorized into eight sub-fields.
 - Language: the language in which the question text is written.
 - Difficulty: a measure of the question's complexity based on the level of reasoning required.

B EXPERIMENT DETAILS

B.1 HUMAN PERFORMANCE ESTIMATION

751Based on the analysis of CFA exams (see https://300hours.com/cfa-passing-score/), the752passing score is approximately 68% for all the three levels. During the annotation process, the753difficulty levels of FAMMA's questions—easy, medium, and hard—closely correspond to those of754CFA Levels I, II, and III. For the medium questions, we assume those who fails the Level I will have755a score of 50% on Level II, therefore the corresponding score for Level II over the whole populationrssis 68% * 68% + (1 - 68%) * 50% = 62.24%. By similarly assuming those who are not qualified for

756 757	STATISTICS	NUMBERS (PERCENTAGE)
758	TOTAL OUESTIONS	1758
759	* MC	976 (55.5%)
760	* OPEN	$782 (44 \ 5\%)$
761	* W EXPLANATIONS	1750 (99.5%)
762	* W MULTIPLE IMAGES	521 (29.6%)
763	TOTAL SUBFIELDS	8
764	* AITEDNATIVE INVESTMENTS	87 (4 9%)
765	* CORDORATE FINANCE	256(14.6%)
766	* DEDIVATIVES	250(14.0%)
767	* Economics	23(10%)
760	* Fourty	256(14.6%)
709	* EVED INCOME	250(14.070)
771	* EDIANGIAL STATEMENT ANALYSIS	99(3.470) 948(14.1%)
772	* PORTEOLIO MANAGEMENT	520(301%)
773	PORTFOLIO MANAGEMENT	529 (50.170)
774	TOTAL IMAGE TYPES	3
775	* QUESTIONS W. TABLES	$1426\ (81.1\%)$
776	* QUESTIONS W. CHARTS	278~(15.8%)
777	* QUESTIONS W. SCREENSHOTS	54 (3.1%)
778	DIFFICULTIES	608 / 438 / 712
779	(F: M: H)	$34.6\% \cdot 24.9\% \cdot 40.5\%$
780	(1. m. n)	
781	Splits	120 / 1638
782	(VALIDATION:TEST)	6.8%:93.2%
783	Languages	1385 / 253 / 120
784	(ENGLISH: CHINESE: FRENCH)	78.8%:14.4%:6.8%
786		
787	AVG LENGTH IN TOKENS	000.40
788	↑ QUESTIONS	233.43
789	* EXPLANATION	73.95

Table 7: More detailed key statistics of FAMMA.

Level III will have a score of 40% on Level III, the expected score of the whole population for Level III becomes 62.24% * 68% + (1 - 62.24%) * 30% = 57.26%. In this context, we set the human score for easy, medium, and hard questions to be equal to that of Level I, II, III—68%, 62.24%, 57.26%, respectively, which resulting a overall score of 59.86%.

B.2 CASE STUDIES

We present a few sample cases of error response from GPT-40.

- Data misinterpretation: GPT-40 incorrectly reads a figure as 16% instead of the correct value of 18%, likely due to a small section of the figure having low resolution from the data source, which leads to inaccurate calculations (see Listing 11).
- Incomplete context understanding: GPT-40 overlooks the option "one of the options are correct" when generating a response, resulting in an incorrect answer (see Listing 12).
- Numerical inaccuracy: GPT-40 produces incorrect decimal values during a square calculation, leading to an erroneous output (see Listing 13).

SUBFIELD	English	CHINESE	French
ALTERNATIVE INVESTMENTS	27 / 17 / 39	3 / 1 / 0	-
CORPORATE FINANCE	72 / 61 / 50	1/35/1	0/33/3
DERIVATIVES	37 / 24 / 172	9/5/0	0/0/7
ECONOMICS	10/1/21	1 / 0 / 0	-
Equity	85 / 24 / 74	20 / 13 / 7	2/28/3
FINANCIAL STATEMENT ANALYSIS	56/21/18	-	-
Fixed Income	76 / 27 / 82	11 / 22 / 12	7 / 7 / 4
PORTFOLIO MANAGEMENT	157 / 59 / 175	31 / 52 / 29	3 / 8 / 15
Total	520 / 234 / 631	76 / 128 / 49	12 / 76 / 32
(2	34.5%:13.4%:52.1%)	(27.7%:50.6%:21.7%)	(7.5%:59.2%:3

Table 8: Distribution of questions in difficulty across languages and subfields in FAMMA.

• Domain knowledge gaps: GPT-40 misunderstands the nature of the high-touch agency approach in financial markets, confusing its application in exchange-traded derivatives and large trades (see Listing 14).

• Ambiguous answer generation: GPT-40 correctly performs the computation but arrives at an incorrect final result due to ambiguity in answer interpretation (see Listing 15).

These instances are firstly categorized by LM-evaluators (see Listing 10 in Appendix B for the instruction prompt), then validated by human expert based on their knowledge and the golden explanations if available.

B.3 DETAILS ON RAG AND 01-REASONING EXPERIMENTS

RAG setup. We utilize 5 CFA Level III curriculum textbooks—"CFA Level III SchweserNotes, Books 1-5, 2023", which comprehensively cover most of the topics found in FAMMA— as the external knowledge source. The notes are in PDF format, each consisting of 200-300 pages with quizzes at the end of every chapter, though these quizzes are not included in FAMMA. We upload them to GPT-40 via the API for queries.

Bynamic COT setup. The implementation is based on the open source project https://github.com/bklieger-groq/g1, which is originally built on Llama-3.1. We improve the project to be compatible with GPT-40. The process begins with a system prompt that includes instructions to guide a step-by-step reasoning approach. Once the problem is introduced as a user message, an assistant message is pre-loaded to establish a standardized starting point for the GPT-40's response generation.

867 868 870 871 { "question_id": " 872 English_validation_86", 873 "context": "The following data are 874 available relating to the performance 875 of Wildcat Fund and the market 876 portfolio: <image_1>", "question": "The risk-free return 877 during the sample period was 7%. >", 878 Calculate Sharpe's measure of 879 performance for Wildcat Fund." 880 "options": "['1.00%', '8.80%', Why?", '44.00%', '50.00%']", "image_1": "/9j/4 AAQSkZJRgABAQAAAQABAAD...]", 883 "image_2": null, "image_3": null, 885 "image_4": null, "image_5": null, "image_6": null, 887 "image_7": null, "image_type": "table", 889 "answers": "C", 890 "explanation": "(18 - 7)/25 = 891 .44.", "topic_difficulty": "easy", 892 "question_type": "multiple-choice", Compressor.", 893 "subfield": "portfolio management", 894 "language": "english", 895 "main_question_id": 369, 896 "sub_question_id": 2, "ans_image_1": null, 897 "ans_image_2": null, "ans_image_3": null 899 Compressor.", 900 Listing 5: Multi-choice questions in JSON 901 representation. 902 903 904 905 906 907 908 909 }, 910 911 912 913

"question_id": " English_validation_42", "context": "Cleveland Compressor and Pnew York Pneumatic are competing manufacturing firms. Their financial statements are printed here .<image_1><image_2><image_3><image_4 "question": "Which firm has the larger investment in current assets? "options": "" "image_1": "/9j/4 AAQSkZJRgABAQAAAQABAAD...", "image_2": "/9j/4 AAQSkZJRgABAQAAAQABAAD...", "image_3": "/9j/4 AAQSkZJRgABAQAAAQABAAD...", "image_4": "/9j/4 AAQSkZJRgABAQAAAQABAAD...", "image_5": null, "image_6": null, "image_7": null, "image_type": "table", "answers": "Cleveland "explanation": "Cleveland Compressor holds the larger investment in current assets. It has current assets of \$92,616 while Pnew York Pneumatic has \$70,101 in current assets . The main reason for the difference is the larger sales of Cleveland "topic_difficulty": "hard", "question_type": "open question "subfield": "financial statement analysis", "language": "english", "main_question_id": 329, "sub_question_id": 3, "ans_image_1": null, "ans_image_2": null, "ans_image_3": null

Listing 6: Open questions in JSON representation.

```
919
         You are a highly knowledgeable
                                                     You are a highly knowledgeable
920
         financial expert. Please answer
                                                     financial expert. Please answer open-
         multiple-choice questions in the
                                                     ended questions in the finance domain.
921
         finance domain. You are given context,
                                                     The questions are multilingual (either
922
         images, questions and options.
                                                     in English, Chinese, or French) and
923
         The questions are multilingual (either
                                                     multimodal (containing images as part
924
                                                     of the question). '<image_1>, <image_2
         in English, Chinese, or French) and
925
         multimodal (containing images as part
                                                     > ...' mentioned in the text of the
         of the question). '<image_1>, <image_2
                                                     context or question are sequential
926
         > ...' mentioned in the text of the
                                                     placeholders for images, which are fed
927
         context or question are sequential
                                                     at the same time as the textual
928
         placeholders for images, which are fed
                                                     information.
929
         at the same time as the textual
                                                     If no image information is provided,
930
         information.
                                                     you must answer based solely on the
                                                     given information.
         If no image information is provided,
931
         you must answer based solely on the
                                                     Besides, the question may contain
932
         given information.
                                                     several sub-questions that share the
933
         Besides, the question may contain
                                                     same context, and the answer to each
934
                                                     sub-question may depend on the answers
         several sub-questions that share the
935
         same context, and the answer to each
                                                     to previous ones.
         sub-question may depend on the answers
                                                     The question format is
936
         to previous ones.
937
         The question format is
                                                     context: <context>
938
                                                     sub-question-1: <sub-question-1>
939
         context: <context>
                                                     sub-question-2: <sub-question-2>
940
         sub-question-1: <sub-question-1>
                                                     sub-question-3: <sub-question-3>
         sub-question-2: <sub-question-2>
941
                                                     . . .
         sub-question-3: <sub-question-3>
942
                                                     Now consider the following question:
943
                                                     context: {context}
944
         Now consider the following question:
                                                     {sub_questions}
945
         context: {context}
                                                     Please provide the answer and a precise
         {sub_questions}
946
                                                     , detailed explanation. The explanation
947
         Please provide the chosen answer and a
                                                      should be in the same language as the
948
                                                     question and should not exceed 400
         precise, detailed explanation of why
949
         this choice is correct. The explanation
                                                     words.
950
          should be in the same language as the
                                                     Your answer must be in a standard JSON
         question and should not exceed 400
                                                     format:
951
         words.
                                                     {{
952
         Your response must be in a standard
                                                        sub-question-1: {{
953
         JSON format:
                                                            answer-1: "answer-1",
954
                                                            explanation-1: "explanation-1"
         {{
955
            sub-question-1: {{
                                                        }},
                answer-1: <answer-1>,
                                                        sub-question-2: {{
956
                explanation-1: <explanation-1>
                                                            answer-2: "answer-2",
957
                                                            explanation-2: "explanation-2"
            }},
958
            sub-question-2: {{
                                                        }},
959
                                                        sub-question-3: {{
                answer-2: <answer-2>,
960
                explanation-2: <explanation-2>
                                                            answer-3: "answer-3",
            }},
                                                            explanation-3: "explanation-3"
961
            sub-question-3: {{
                                                        }},
962
                answer-3: <answer-3>,
963
                                                     }}
                explanation-3: <explanation-3>
964
            }},
                                                     Ensure that the response strictly
                                                     adheres to JSON syntax without any
965
         }}
                                                     additional content.
966
         Ensure that the response strictly
967
                                                    Listing 8: Format of our instruction prompt on open
         adheres to JSON syntax without any
968
                                                    questions.
         additional content.
969
```

970 Listing 7: Format of our instruction prompt on multi-

971 choice questions.

973 974 975 976 977 You are a highly knowledgeable expert You are a highly skilled expert in 978 and teacher in the finance domain. error analysis for AI models in the 979 finance domain. You are reviewing You are reviewing a student's answers 980 to financial questions. collected incorrect answers to 981 The questions are multilingual (either financial questions. 982 in English, Chinese, or French) and The questions are multilingual (either multimodal (containing images as part in English, Chinese, or French) and 983 of the question). '<image_1>, <image_2 multimodal (containing images as part 984 > ...' mentioned in the text of the of the question). '<image_1>, <image_2 985 context or question are sequential > ...' mentioned in the text of the 986 placeholders for images, which are fed context or question are sequential at the same time as the textual placeholders for images, which are fed 987 information. at the same time as the textual 988 You are given the context, the question information. 989 You are given the context, the question , the student's answer and the student' 990 s explanation and the ground-truth , the student's answer, the student's 991 explanation and the ground-truth. answer. Please use the given information and 992 refer to the ground-truth answer to You need to classify these incorrect 993 determine if the student's answer is answers based on the provided 994 correct. categories: perceptual errors, lack of 995 knowledge, reasoning errors, and other 996 The input information is as follows: errors. Here are the definitions for 997 each error type: context: {context} 998 question: {question} Data misinterpretation: .(omitted) 999 student's answer: {model_answer} Incomplete context understanding: ...(1000 student's explanation: { omitted) 1001 model_explanation} Numerical inaccuracy: ...(omitted) Domain knowledge gaps: ...(omitted) ground-truth answer: {answer} 1002 Ambiguous answer generation: ...(1003 Please respond directly as either ' omitted) 1004 correct' or 'incorrect'. 1005 The input is as follows; use these Listing 9: Format of our prompt on judging the details to determine the primary error correctness of the model output. 1007 category. 1008 context: {context} 1009 question: {question} 1010 model incorrect answer: {model_answer} 1011 model explanation: {model_explanation} ground-truth answer: {answer} 1012 1013 Now please provide the result directly, 1014 identifying the error category as one 1015 of: data misinterpretation, incomplete 1016 context understanding, numerical inaccuracy, domain knowledge gaps, or 1017 ambiguous answer generation. 1018 1019 Listing 10: Format of our prompt on error analysis on 1020 model's output. 1021 1023 1024 1025

```
1026
1027
1028
1029
1030
1031
          The expected return (E[R]) is
                                                          Therefore, the expected holding-period
1032
          calculated as follows:
                                                          return for the stock is 10.4%. However,
1033
                                                          since 10.4\% is not one of the given
         E[R] = (0.30 * 16\%) + (0.50 * 12\%) +
1034
          (0.20 * -5\%) = 4.8\% + 6\% - 1\% = 9.8\%
                                                           options, the closest correct answer is
1035
                                                           8.33\%, which is option B.
        Listing 11: A sample error response of GPT-40 due to
1036
        data misinterpretation: it incorrectly reads the figure as Listing 12: A sample error response of GPT-40 due
1037
                                                         to incomplete context understanding: it overlooks the
        16\% instead of the correct value of 18\%.
                                                         option of "None of the options are correct".
1039
1040
1041
1042
1043
1045
1046
1047
1048
1049
         For Normal growth: 0.45 * (15% -
                                                          Small currency trades and small
1050
         14.35\%)^2 = 0.45 * (0.65\%)^2 = 0.45 *
                                                           exchange-traded derivatives trades are
         0.004225 = 0.00190125
                                                           typically implemented using the direct
1051
                                                          market access (DMA) approach, and the
1052
        Listing 13: A sample response of GPT-40 due to
                                                          high-touch agency approach is typically
1053
        numerical inaccuracy: the square calculation results in
                                                          used to execute large, non-urgent
1054
        incorrect decimal values.
                                                           trades in fixed-income and
                                                           exchange-traded derivatives markets.
1055
1056
                                                         Listing 14: A sample case response of GPT-40
1057
                                                         due to domain knowledge gaps: the high-touch
1058
                                                         agency approach is mainly used for illiquid orders but
                                                         exchange-traded derivatives tend to be very liquid.
1062
1063
1064
1065
1067
1068
1069
          The calculation is as follows: (0.30 \ast
1070
          18\%) + (0.50 * 12\%) + (0.20 * -5\%) =
1071
         5.4\% + 6\% - 1\% = 10.4\%. Therefore, the
1072
         expected HPR for KMP stock is 10.88%.
        Listing 15: A sample case response of GPT-40 due to
1074
        ambiguous answer generation: it correctly performs the
1075
        computation but reaches to a wrong final result.
1076
1077
1078
```

1080 1081 1082 1083 1084 1085 1086 1087 1088 { 1089 "role": "system", "role": "assistant", "content": """You are an expert AI "content": "Thank you! I will now 1090 assistant who explains your reasoning think step by step following my 1091 process step by step. For each step, instructions, starting at the beginning 1092 after decomposing the problem." provide a title describing what you're 1093 doing in that step, and the content. ł 1094 Determine whether another step is Listing 17: Format of the assistant prompt used in needed or if you're ready to give a 1095 dynamic COT. final answer. Respond in JSON format 1096 with 'title', 'content', and ' next_action' (which can be 'continue' { 1098 or 'final_answer') keys. Use multiple "role": "user", 1099 reasoning steps whenever possible, at "content": "Please provide the least 3. Be aware of your limitations 1100 final answer based on the above as an LLM and what you can and cannot 1101 reasoning.", do. In your reasoning, include 1102 exploration of alternative answers. 1103 Consider that you might be wrong and Listing 18: Format of the user prompt used in dynamic 1104 where errors in your reasoning might COT occur. Thoroughly test all other 1105 possibilities. You may be wrong. When 1106 you say you're revisiting, actually 1107 revisit and use a different method to 1108 do so. Don't just say you're revisiting 1109 . Use at least 3 methods to arrive at an answer. Use best practices. 1110 1111 Example of a valid JSON response: 1112 '''json 1113 { "title": "Identifying Key 1114 Information". 1115 "content": "To begin solving 1116 this problem, we need to carefully 1117 examine the given information and 1118 identify the key elements that will guide our solution process. This 1119 involves...", 1120 "next_action": "continue" 1121 } . . . 1122 1123 } 1124 Listing 16: Format of the system prompt used in 1125 dynamic COT. 1126 1127 1128 1129 1130 1131 1132 1133