

Learning to Trade Like an Expert: Cognitive Fine-Tuning for Stable Financial Reasoning in Language Models

Yuchen Pan
ypan@link.cuhk.edu.hk
Department of Information Engineering
The Chinese University of Hong Kong
Hong Kong SAR, China

Soung Chang Liew
soug@ie.cuhk.edu.hk
Department of Information Engineering
The Chinese University of Hong Kong
Hong Kong SAR, China

Abstract

Recent deployments of large language models (LLMs) as autonomous trading agents raise questions about whether financial decision-making competence generalizes beyond specific market patterns and how it should be trained and evaluated in noisy markets lacking ground truth. We propose a structured framework for training and evaluating such models. Central to our approach is a curated, multiple-choice question (MCQ) dataset derived from classic textbooks and historical markets, verified by an AI committee, enriched with structured reasoning traces, and augmented to reduce shortcut learning. To evaluate whether performance on isolated MCQs generalizes to real-world trading, we introduce a two-stage protocol combining test-set evaluation with an MCQ-based chronological trading simulation. Extensive evaluations across market regimes provide statistically robust evidence that open models trained with our framework exhibit competitive, risk-aware behavior over time, outperform open-source baselines, and approach frontier-model performance at smaller scale. We release the dataset and evaluation framework to support further research.

1. Introduction

In late October 2025, financial artificial intelligence reached a new milestone with the launch of the Alpha Arena initiative (<https://nofl.ai/>). For the first time, frontier LLMs were entrusted with real capital to trade autonomously in live U.S. equity and cryptocurrency markets, competing directly against one another.

Despite the widespread attention it received, the initiative raises fundamental questions about autonomous decision-making in finance. **Question 1:** Does the observed profitability reflect genuine domain competence that generalizes across diverse market conditions, or is it merely an artifact of specific market patterns tied to particular time periods? **Question 2:** If such competence exists, can it be achieved more reliably through specialized trading models rather than general-purpose LLMs? Even if these agents demonstrate generalizable trading competence, deploying them in institutional finance presents additional challenges related to training, governance, and control. **Challenge 1:** Reliance on externally hosted, closed-source models introduces significant risks, including data privacy concerns, regulatory compliance issues, and the potential leakage of proprietary trading strategies—factors that increasingly motivate a shift toward locally deployed models. **Challenge 2:** Financial markets are inherently noisy and lack clear “correct” answers, making supervised learning on raw price data prone to overfitting spurious correlations rather than capturing causal decision principles. As a result, financial institutions have a compelling need for alternative training methods and deployment platforms that provide high-quality guidance while fostering reasoning and judgment under uncertainty.

Existing literature offers few direct precedents to address these questions. Most open-source financial language models are not designed for autonomous execution: Discriminative models (Araci, 2019; Huang et al., 2023; Liu et al., 2021) primarily focus on sentiment analysis, while generative models (Wu et al., 2023; Yang et al., 2023a; Zhang et al., 2023) are largely aimed at semantic retrieval or advisory tasks (Yang et al., 2023b; Chen et al., 2023). Recent execution-oriented frameworks (Zhou et al., 2024; Zhang et al., 2024; Xiao et al., 2025) often rely on complex multi-agent systems,

This work was supported in part by a CUHK Direct Grant (4055243).
Accepted to FoGen 2026: Foundations of Deep Generative Models: Understanding Memorization, Generalization, and Reasoning, an ICML 2026 workshop (non-archival).

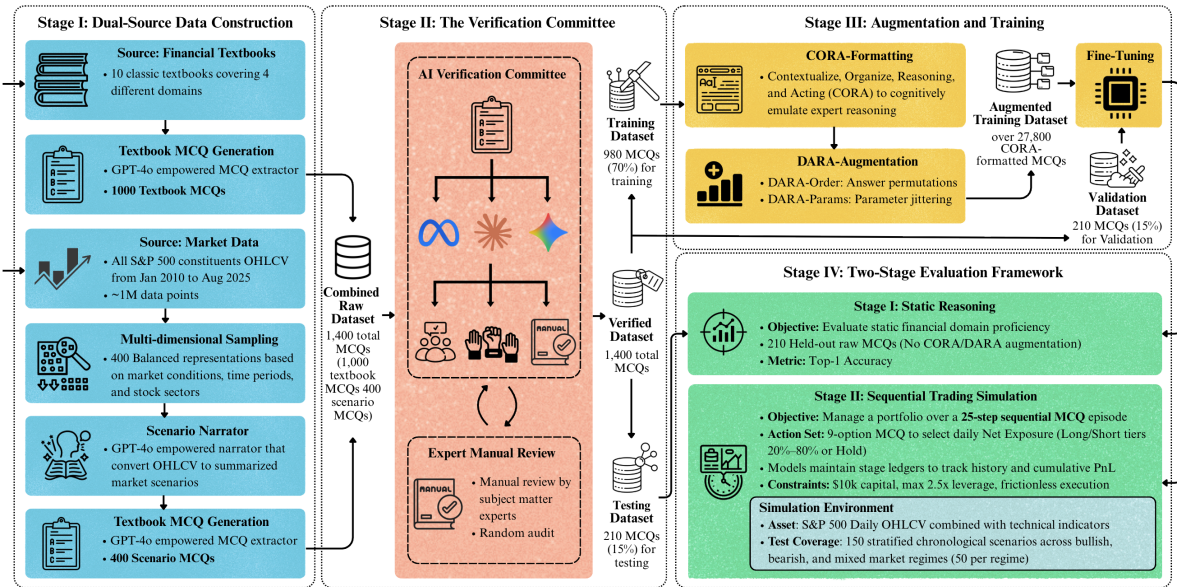


Figure 1. Overview of the proposed framework.

which obscure the model’s intrinsic reasoning process. Furthermore, existing financial datasets and benchmarks (Chen et al., 2021; Li et al., 2024) predominantly emphasize data extraction or price prediction rather than decision-making logic. To our knowledge, no publicly available resource explicitly provides a framework for financial decision logic.

Building on these questions and challenges, we make three contributions toward training and evaluating robust, locally deployable financial agents. **First**, in response to Question 2 and Challenge 2, we introduce a curated MCQ dataset for training and evaluating autonomous trading agents. The dataset is derived from classic textbooks and historical market scenarios grouped by market regime, time period, and sector. Each item is verified by an AI committee composed of three LLMs for factual and logical errors, annotated with a four-step chain-of-thought (CoT) trace, and augmented by shuffling answer options and varying technical indicator parameters to reduce shortcut learning. **Second**, to investigate Question 1, we propose a two-stage evaluation protocol: The first stage measures performance on isolated MCQs in a held-out test set, while the second stage evaluates whether this competence generalizes to better real-world trading behavior through a chronological, MCQ-based trading simulation under realistic portfolio constraints. Extensive evaluations across bullish, bearish, and mixed market regimes provide statistically robust evidence that models trained with our framework translate strong first-stage performance into consistent, risk-aware real-world trading behavior, enabling smaller-scale open models to perform competitively with controlled downside risk. **Finally**, to bridge the deployment gap in Challenge 1, we open source the full training and evaluation pipeline at <https://github.com/Lauchlan-Pan/CORA-Learning-to-trade-like-an-expert> to support customization and further research on locally deployed financial agents. Figure 1 summarizes the proposed framework.

2. Related Works

Financial LLM Evolution Financial language models have progressed from task-specific classifiers to more general systems used in financial analysis and trading applications (Li et al., 2023; Ding et al., 2024; Dong et al., 2025). Early models such as FinBERT (Araci, 2019; Liu et al., 2021; Huang et al., 2023) were primarily designed for sentiment analysis, mapping financial news and social media to trading signals. For broader financial applications beyond sentiment analysis, closed-source models like BloombergGPT (Wu et al., 2023) were pre-trained for tasks such as financial question answering and investment advisory. Open-source models including FinGPT (Yang et al., 2023a; Zhang et al., 2023), InvestLM (Yang et al., 2023b), and DISC-FinLLM (Chen et al., 2023) have further adapted open models to these tasks through parameter-efficient post-training. More recently, LLMs have been integrated directly into trading pipelines via hierarchical planning (Zhou et al., 2024) and multi-agent designs (Zhang et al., 2024; Yu et al., 2024; Xiao et al., 2025). In this work, we combine open-weight post-training with direct trading execution, relying on a single model instead of complex planning or multi-agent architectures.

Structured Reasoning Strategies Another line of work focuses on improving the reliability of LLM outputs, which is especially important in finance. Prior studies showed that standard LLMs can produce hallucinated explanations that are not causally aligned with their final answers (Webb et al., 2023; Binz & Schulz, 2023; Turpin et al., 2023). To mitigate this, inference-time methods prompt models to produce explicit intermediate steps, including linear reasoning chains (Wei et al., 2022; Wang et al., 2022), branching sets of reasoning paths (Yao et al., 2023; Besta et al., 2024), and iterative self-critique and revision (Shinn et al., 2023; Madaan et al., 2023). Training-time methods encourage similar behavior by fine-tuning models on annotated reasoning traces (Uesato et al., 2022; Hsieh et al., 2023; Hao et al., 2023; Lightman et al., 2023) or by sampling multiple candidate outputs and fine-tuning the model on the highest-quality ones (Zelikman et al., 2022; 2024). Our work follows the training-time paradigm by fine-tuning models on annotated reasoning traces.

Data Curation and Benchmarks Training and evaluating financial models is challenging due to market noise and the lack of definitive ground truth. Prior work showed that problems with clear structure and explicit intermediate steps can yield reliable signals for learning and assessment (Gunasekar et al., 2023; Wang et al., 2025). This principle appears in financial datasets like FinQA and ConvFinQA (Chen et al., 2021; 2022), which emphasize financial text and numerical tables, as well as MultiHiertt and BizBench (Zhao et al., 2022; Krumdick et al., 2024), which focus on longer business documents. In parallel, datasets based on real market data—including AlphaFin (Li et al., 2024), FinBen (Xie et al., 2024), and StockBench (Chen et al., 2025)—emphasize time-series signals and typically frame tasks as information retrieval or market prediction. Our work adopts both directions by deriving problems from real market data and annotating them with explicit reasoning steps.

3. Cognitive Financial Reasoning Dataset

3.1. Dual-Source Data Construction

To support structured training and evaluation of trading agents, we constructed the Cognitive Financial Reasoning Dataset. We used an MCQ format to define a clear decision space that supports controlled analysis of model choices. As illustrated in Figure 1 (Stage 1), our data construction pipeline is designed to mimic the *learning process of human experts*: first acquiring foundational principles from classic **Textbooks**, and then refining judgment through exposure to dynamic **Market Data**. The full list of source textbooks, the prompt templates used to generate MCQs from textbooks and market data, and the final datasets are available in our open-source repository.

3.1.1. TEXTBOOK KNOWLEDGE EXTRACTION

To build a reliable foundation for financial reasoning, we first curated ten classic textbooks spanning four key domains essential to professional trading practice: Technical Analysis, Quantitative Trading, Macroeconomics, and Trading Psychology. We used GPT-4o (Hurst et al., 2024) to identify core concepts within each textbook section and then converted them into discrete MCQs that emphasize reasoning about trading decisions rather than simple recall of definitions or facts. For example, for a chapter introducing indicator divergence, we generated an MCQ describing a price downtrend with a bullish RSI divergence and asked for the most prudent trading action. We generated exactly 100 MCQs per textbook, yielding a balanced theoretical core of 1,000 textbook-derived questions.

3.1.2. MARKET SCENARIO EXTRACTION

To help models make decisions in noisy, real-world markets, we created 400 additional scenario-based MCQs using roughly 15 years of S&P 500 historical data (2010–2025). We started with about one million data records and split the time series into fixed-length windows. Each window was annotated with three tag categories (listed below), and we sampled roughly evenly across the resulting tag tuples to ensure balanced coverage across:

1. Market regimes (e.g., strong uptrends, high volatility)
2. Time periods (e.g., post-crisis recovery)
3. Market sectors (e.g., healthcare, energy)

We used a two-step GPT-4o pipeline to convert numeric signals into text. For each price window, the **Narrator** summarized the data and 38 technical indicators into a natural-language description (e.g., “*TSLA shows overbought*”).

conditions with $%K=97.6$; a slightly negative $MACD=-0.12$ suggests waning bullish momentum. . .”). The **Generator** then turned the summary into an MCQ asking for the best trading action (e.g., enter now). Together with the textbook questions, this produced a pre-verification dataset of 1,400 MCQs.

3.2. The Verification Committee

To mitigate hallucinations and over-reliance on the single generator model GPT-4o, we introduce a verification stage (Figure 1, Stage 2) using an **AI committee** of three prior-generation frontier LLMs: Gemini-2.5-Flash (Comanici et al., 2025), Claude-3.7-Sonnet (Anthropic et al., 2025b), and Llama-4-Scout (Meta, 2025). This setup is similar to using multiple human judges to reduce individual errors. We use earlier-generation models because our later experiments evaluate newer frontier models as baselines; keeping the committee separate helps avoid overlap between verification and evaluation. To safeguard label quality, expert human review is used as the final source of labels for difficult cases.

Each of the 1,400 raw MCQs was independently evaluated by all AI committee members under a standardized prompt (released in our repository). For every item, each model was required to choose an answer, provide a brief rationale, and flag any suspected factual errors; when the information was insufficient, it was instructed to output **UNCERTAIN** rather than guess. Final labels were assigned via a strict consensus rule:

- **Unanimous agreement (3/3):** 1,139 items (81.4%) were retained as high-confidence labels.
- **Majority agreement (2/3):** 158 items (11.3%) were accepted, with the committee’s majority answer replacing the original generator’s label.
- **No agreement:** The remaining 103 items were reviewed by domain experts; 96 reached consensus after revising the MCQs to clarify wording and correct factual issues, and 7 were further simplified for clarity.

In a follow-up random audit of 30 verified items, human experts agreed with the committee’s labels on 29 of them (96.67%), suggesting that the AI committee can approximate expert verification at scale.

3.3. CORA: Cognitive Chain-of-Thought

For humans, expertise in financial decision-making is not acquired by memorizing isolated answers, but through a progressive learning process of first reading market context, then framing the decision, weighing alternatives, and finally translating judgment into a concrete trading plan. To mirror this learning trajectory, we introduce **CORA (Contextualize–Organize–Reason–Act)**, a four-stage CoT template that expands each verified MCQ into:

- **Contextualize:** Synthesizes the market state into a coherent representation, integrating trends, volatility, technical indicators, and relevant background factors.
- **Organize:** Defines the decision frame by specifying the trading objective and key constraints, reflecting how experts simplify complex environments.
- **Reason:** Justifies the chosen option and explicitly explains why the other options are less appropriate, encouraging contrastive reasoning.
- **Act:** Translates the conclusion into a concrete trading plan, including direction, entry conditions, and risk controls.

Figure 2 illustrates a simplified example of how a single MCQ is expanded into CORA format. We use GPT-4o with standardized prompts to generate the four-stage traces. Manual audits of 50 randomly selected items confirm that the generated reasoning is clear and consistent with the underlying question. By turning simple answer labels into structured reasoning traces, CORA shifts the learning signal from *which* decision is correct to *how* decisions are formed under uncertainty. Instead of treating financial markets as if they had definitive ground truth, CORA encodes expert-like patterns—such as attention to volatility and explicit risk management—that resemble how human traders reason in practice.

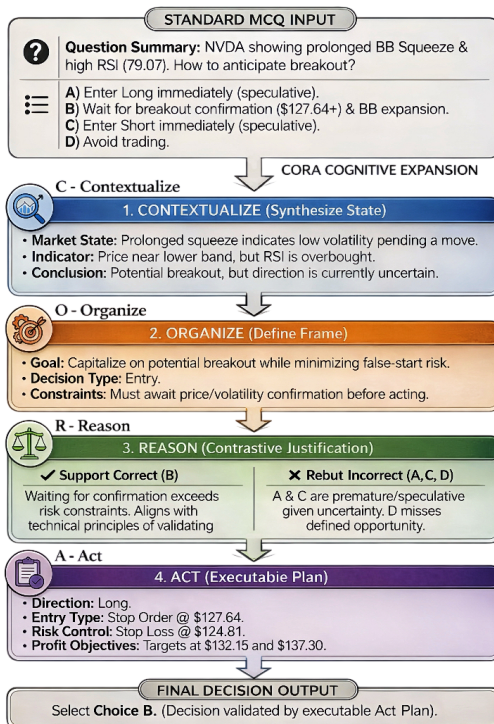


Figure 2. A simplified example of CORA cognitive expansion.

3.4. DARA: Dual-Axis Robustness Augmentation

Analysis of the initial dataset revealed two systematic biases introduced during generation. First, correct answers were not evenly spread across options: In the 1,400 MCQs, options A and B covered 76.9% of correct labels, while option D appeared in only 4.2% of cases. Second, many textbook-derived questions reused the same indicator settings (e.g., repeatedly using a 14-day moving average). These patterns can lead to shortcut learning, where models exploit superficial cues such as “*A is usually correct*” or fixed parameter values instead of reasoning about the underlying scenario.

To address these issues, we introduce **Dual-Axis Robustness Augmentation (DARA)**, which preserves the CORA structure while reducing reliance on answer position and specific numeric settings. It operates along two dimensions:

- **DARA-Order:** Addresses positional bias by generating all $4! = 24$ permutations of the answer choices for each MCQ.
- **DARA-Params:** Reduces overfitting to specific indicator values by introducing controlled variations within common practical ranges (e.g., varying RSI-14 to RSI-21).

When applied to the 980 training MCQs, DARA-Order yields a $24\times$ expansion, while DARA-Params introduces an additional $\sim 1.18\times$ increase, resulting in approximately 27,800 training examples. We do not apply CORA or DARA to the validation or test sets; they remain unchanged to provide a consistent benchmark for measuring generalization.

Together, CORA and DARA form our reasoning-focused augmentation strategy: CORA makes the decision logic explicit, while DARA discourages reliance on answer position and fixed numeric settings. The prompt templates for both stages are included in our open-source repository.

4. Two-Stage Evaluation Framework

We propose a two-stage evaluation framework for financial reasoning that evaluates models on both **static reasoning**, where each decision is an independent MCQ response, and **sequential trading behavior**, where decisions accumulate

over time in a live portfolio. At test time, models receive only the MCQ prompt (market scenario and answer choices in Stage I, plus the current portfolio ledger state in Stage II), and no gold CORA traces are provided. Models are prompted to generate their own reasoning (CORA-style for CORA-trained variants; standard CoT for others) before giving a final answer.

4.1. Stage I: Static Reasoning

Stage I follows a standard supervised evaluation paradigm. We evaluate each model on a test set of 210 MCQs. Each question is answered independently to assess how well the model performs on isolated financial decisions.

4.2. Stage II: Sequential Trading Simulation

While Stage I evaluates isolated MCQ answering accuracy, static performance alone is insufficient for assessing practical trading competence. In real financial markets, decisions unfold sequentially, affecting future capital, exposure, and risk. Models that perform well on static tasks may nevertheless behave unstably or inconsistently when their decisions are chained together in a portfolio.

To address this, Stage II introduces a **chronological evaluation protocol** that preserves the MCQ format of Stage I while placing decisions within a continuous trading simulation. In Stage II, the model is asked to manage a portfolio over a fixed 25-step episode (approximately one trading month). At every step, it receives the current market scenario along with a state ledger documenting prior actions, current positions, and cumulative profit and loss. This running context allows each decision to condition on earlier outcomes. The model then answers an MCQ by selecting one of nine target net exposures for the next session: long or short at 20%, 40%, 60%, or 80% of allowable exposure, or a neutral hold.

Stage II also differs from open-ended agent setups (such as the Alpha Arena initiative described in the Introduction), where models operate in large, unstructured action spaces using free-form text, making it hard to separate reasoning quality from parsing errors or prompt-specific artifacts. In contrast, Stage II restricts the action space to a small set of interpretable exposure choices, making it easier to relate each decision to a clear change in the portfolio.

Episodes are built from daily S&P 500 data and standard technical indicators using the same pipeline as the scenario-based MCQs, except that the answer choices are always fixed to the same nine exposure options mentioned above. Models receive no access to future information. We enforce an initial capital of \$10,000, a maximum gross leverage of $2.5\times$, and frictionless execution to isolate decision quality from trading frictions such as bid-ask spreads and slippage.

To enable **statistically robust evaluation** beyond specific market patterns, we run 150 non-overlapping episodes across bullish, bearish, and mixed market regimes (50 per regime), for a total of 3,750 MCQ steps. Regimes are defined by buy-and-hold returns over each 25-day window, and episodes are selected using a round-robin strategy over time to avoid overrepresenting any particular period.

5. Experiments and Results

5.1. Experimental Setup

We fine-tune Llama-3.1-8B-Instruct (Meta, 2024) on our Cognitive Financial Reasoning Dataset using Q-LoRA (Dettmers et al., 2023). The model is selected for its open-weight availability and strong instruction-following performance at the 8B scale, while Q-LoRA enables efficient fine-tuning under resource-constrained settings. Training uses the standard Llama-3.1 chat template: "system messages" in the template define the task, "user messages" provide market scenarios and MCQ options, and "assistant messages" output a CORA-structured reasoning trace and the final selected option (i.e., the target completion). Each MCQ is serialized as a single chat turn, and cross-entropy is applied only to assistant tokens (system/user tokens are masked), encouraging learned reasoning and answer selection rather than prompt reproduction.

We benchmark against frontier LLMs and size-matched open-source models. Frontier LLMs (e.g., GPT-5.1 (OpenAI et al., 2025)), similar to those used in the Alpha Arena initiative, are evaluated via APIs and serve as reference points for high-capacity models on financial tasks. Open-source models at comparable scales (e.g., Qwen3-8B (Yang et al., 2025)) serve as size-matched baselines for locally deployable systems. For both Stage I and Stage II (including all

Table 1. Stage I and Stage II performance. Stage I reports test accuracy. Stage II reports overall average return (Ret, %), mean net exposure (Exp; average absolute fraction of capital allocated), and turnover (Tur; cumulative absolute change in net exposure per episode, %), together with regime-specific average returns and downside risk measured by CVaR-5% (%), defined as the mean return of the worst-performing 5% episodes. Stage II results are not reported for *w/o CORA* due to degenerate behavior (the model collapsed into selecting identical exposure levels).

Model / Method	Stage I	Stage II: Overall				Stage II: Regime-Specific					
	Static	(150 Episodes)				Bullish		Bearish		Mixed	
	Acc	Ret	Exp	Tur	Ret	CVaR	Ret	CVaR	Ret	CVaR	
Frontier Models											
GPT-5.1 (OpenAI et al., 2025)	80.95	7.08	0.68	157.7	10.52	-13.1	9.81	-3.41	0.91	-7.89	
Claude-4.5-Sonnet (Anthropic et al., 2025a)	86.19	12.04	0.82	54.5	21.26	3.80	14.59	-0.24	0.26	-9.60	
DeepSeek-v3.1 (DeepSeek-AI et al., 2024)	87.14	10.56	1.30	64.8	45.84	3.75	-14.45	-33.43	0.29	-4.23	
Open-Source Baselines											
Llama-3.1-8B-Instruct	57.62	6.23	1.11	284.6	24.61	-14.35	-4.93	-32.60	-0.99	-14.68	
Qwen3-8B (Yang et al., 2025)	71.90	6.33	0.87	185.2	10.71	-21.42	8.05	-25.23	0.24	-13.96	
Mistral-7B-Instruct-v0.3 (Jiang et al., 2023)	64.29	0.43	1.19	315.2	20.12	-11.17	-19.41	-47.57	0.57	-12.99	
Ablation Analysis											
Ours w/o CORA, DARA & Verification	70.95	1.36	1.22	141.2	27.57	1.04	-14.25	-28.36	0.19	-5.42	
Ours w/o DARA	76.67	3.07	0.96	29.5	22.39	3.81	-13.39	-29.37	0.22	-1.69	
Ours w/o CORA	46.19	—	—	—	—	—	—	—	—	—	
Ours (Full Model)	82.38	7.64	1.70	82.7	43.66	8.19	-21.67	-36.94	0.93	-2.64	

150 trading episodes), we set the decoding temperature to 0 for all baseline and locally fine-tuned models to ensure deterministic and directly comparable outputs. To isolate the effects of the components in our method, we run ablations that remove verification, CORA, and DARA. For a fair comparison, CORA-trained models are prompted to produce the full four-stage CORA output before giving a final answer, whereas other models are prompted to generate standard CoT with a similar output length. This setup ensures that performance differences primarily reflect the learned reasoning structure, rather than variations in the amount of intermediate text generated.

5.2. Results and Discussion

Table 1 reports performance for both Stage I and Stage II.

Stage I Frontier models achieve high MCQ accuracy (80.95–87.14%), while open-source baselines lag behind. Our full model attains 82.38% accuracy, outperforming all open-source baselines under the same conditions and approaching the frontier range despite its smaller size and lack of external tools. Ablation results clarify the sources of these gains. Standard supervised fine-tuning on raw MCQs (*without CORA, DARA and verification*) achieves 70.95% accuracy, whereas the full pipeline improves this by over 11 percentage points. Removing DARA lowers accuracy to 76.67% and increases shortcut behavior. Removing CORA has a much stronger effect: Accuracy drops to 46.19%, and the model overfits to the DARA permutations, collapsing into repeatedly selecting the same option. These results suggest that our structured training is crucial when learning from noisy financial data.

Stage II We next test whether single-step performance carries over to sequential trading using metrics that capture profitability, capital use, behavioral stability, and downside risk (defined in the table caption). In the overall evaluation (150 episodes), our full model achieves an average return of 7.64%, outperforming all open-source baselines and approaching frontier-level performance. It also maintains the highest mean net exposure with moderate turnover, indicating confident positioning without excessive trading. Several baselines either trade too frequently or keep exposure too low, suggesting difficulty in maintaining stable policies. Results for the *w/o CORA* variant are omitted because it exhibits degenerate behavior, selecting identical exposure levels across all episodes (consistent with the Stage I findings).

Regime Sensitivity Breaking the results down by market regime shows how model behavior changes across conditions. In bullish regimes, our model delivers strong returns (43.66%) and achieves the best 5% CVaR, indicating controlled downside risk. In mixed regimes, it again attains the highest average return with limited tail risk, suggesting it can

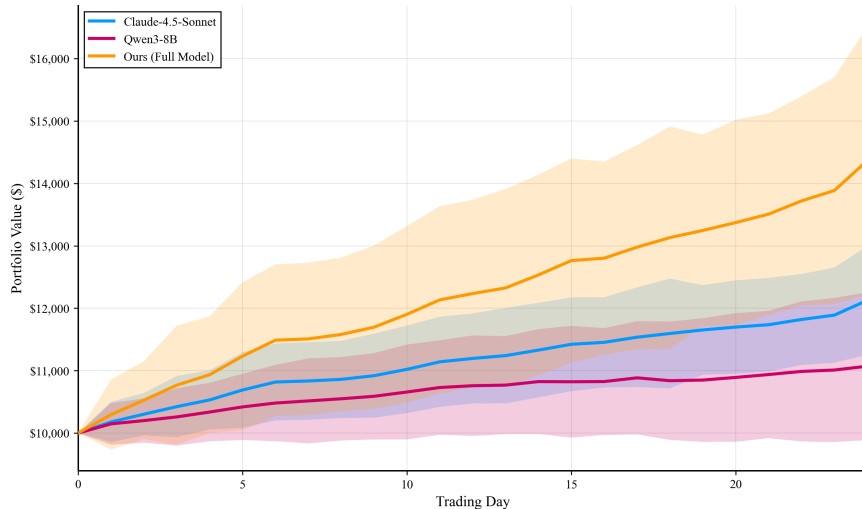


Figure 3. **Bullish-regime wealth trajectories.** Solid lines show mean portfolio value over 50 simulation runs for Claude-4.5-Sonnet, Qwen3-8B, and our Full Model; shaded regions show ± 1 standard deviation.

handle uncertainty without taking excessive downside risk. In bearish regimes, however, it underperforms most baselines in both average return and tail risk. Overall, these results suggest an assertive long-biased policy—also seen in DeepSeek-v3.1—that works well in favorable or mixed conditions but is not well adapted to sustained bearish markets. Figure 3 visualizes portfolio value trajectories in bullish-regime episodes. For clarity, we plot three representative models with the highest overall return in each category. Shaded regions indicate dispersion across 50 runs within ± 1 standard deviation. The Full Model exhibits a consistently higher distribution of outcomes across the trading horizon, suggesting that the observed advantage reflects a systematic behavioral improvement rather than a purely random fluctuation.

Stage I vs. Stage II Relationship Stage I accuracy and Stage II trading performance are related but not equivalent. Among frontier models, the highest Stage I accuracy comes from DeepSeek-v3.1 (87.14% accuracy, 10.56% return), but the best average return is achieved by Claude-4.5-Sonnet (86.19%, 12.04%). This shows that stronger single-step accuracy does not always translate into better portfolio outcomes. However, within our variants, improvements in Stage I accuracy do track improvements in Stage II performance. Moving from *standard fine-tuning* (70.95%, 1.36%) to *without DARA* (76.67%, 3.07%), and then to the *full model* (82.38%, 7.64%), increases both test accuracy and average return. The ablations also highlight that removing DARA keeps Stage I accuracy relatively high but reduces trading performance. Overall, high static accuracy appears necessary but not sufficient for strong trading behavior, and the CORA- and DARA-based training pipeline helps translate single-step reasoning into more consistent portfolio decisions over time.

6. Conclusion

This paper investigates whether LLMs can acquire generalizable financial decision-making competence beyond specific market patterns. We introduce the Cognitive Financial Reasoning Dataset, along with an AI verification committee, CORA, a four-stage reasoning template, and DARA, a data augmentation method that reduces shortcut learning in noisy market settings. To evaluate whether single-step reasoning transfers to real-world trading, we propose a two-stage evaluation framework that combines static MCQ testing with a chronological trading simulation under portfolio constraints.

Our results show that models trained with this pipeline outperform open-source baselines and approach frontier-model performance at a smaller scale in bullish and mixed regimes, but require further improvement in bearish regimes.

Future work includes strengthening robustness in bearish markets, extending the framework to multi-asset settings, adding signals such as market sentiment and trading friction, and studying how explicit reasoning traces can support interpretability, governance, and safe real-world deployment.

References

- Anthropic et al., Sep 2025a. URL <https://assets.anthropic.com/m/12f214efcc2f457a/original/Claude-Sonnet-4-5-System-Card.pdf>.
- Anthropic et al. Claude 3.7 sonnet system card, February 2025b. URL <https://www.anthropic.com/claude-3-7-sonnet-system-card>. Released February 23, 2025.
- Araci, D. Finbert: Financial sentiment analysis with pre-trained language models. *arXiv preprint arXiv:1908.10063*, 2019.
- Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Podstawski, M., Gianinazzi, L., Gajda, J., Lehmann, T., Niewiadomski, H., Nyczyk, P., et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pp. 17682–17690, 2024.
- Binz, M. and Schulz, E. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120, 2023.
- Chen, W., Wang, Q., Long, Z., Zhang, X., Lu, Z., Li, B., Wang, S., Xu, J., Bai, X., Huang, X., and Wei, Z. Disc-finllm: A chinese financial large language model based on multiple experts fine-tuning. *arXiv preprint arXiv:2310.15205*, 2023.
- Chen, Y., Yao, Z., Liu, Y., Ye, J., Yu, J., Hou, L., and Li, J. Stockbench: Can LLM Agents Trade Stocks Profitably In Real-world Markets? 2025. doi: 10.48550/ARXIV.2510.02209. URL <https://arxiv.org/abs/2510.02209>.
- Chen, Z., Chen, W., Smiley, C., Shah, S., Borova, I., Langdon, D., Moussa, R., Beane, M., Huang, T.-H., Routledge, B. R., et al. Finqa: A dataset of numerical reasoning over financial data. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 3697–3711, 2021.
- 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. *arXiv preprint arXiv:2210.03849*, 2022.
- Comanici, G., Bieber, E., Schaekermann, M., Pasupat, I., Sachdeva, N., Dhillon, I., Blistein, M., Ram, O., Zhang, D., Rosen, E., et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- DeepSeek-AI, Liu, A., Feng, B., Xue, B., Wang, B., Wu, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., Dai, D., Guo, D., Yang, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H., Bao, H., Xu, H., Wang, H., Zhang, H., Ding, H., Xin, H., Gao, H., Li, H., Qu, H., Cai, J. L., Liang, J., Guo, J., Ni, J., Li, J., Wang, J., Chen, J., Chen, J., Yuan, J., Qiu, J., Li, J., Song, J., Dong, K., Hu, K., Gao, K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Xu, L., Xia, L., Zhao, L., Wang, L., Zhang, L., Li, M., Wang, M., Zhang, M., Zhang, M., Tang, M., Li, M., Tian, N., Huang, P., Wang, P., Zhang, P., Wang, Q., Zhu, Q., Chen, Q., Du, Q., Chen, R. J., Jin, R. L., Ge, R., Zhang, R., Pan, R., Wang, R., Xu, R., Zhang, R., Chen, R., Li, S. S., Lu, S., Zhou, S., Chen, S., Wu, S., Ye, S., Ye, S., Ma, S., Wang, S., Zhou, S., Yu, S., Zhou, S., Pan, S., Wang, T., Yun, T., Pei, T., Sun, T., Xiao, W. L., Zeng, W., Zhao, W., An, W., Liu, W., Liang, W., Gao, W., Yu, W., Zhang, W., Li, X. Q., Jin, X., Wang, X., Bi, X., Liu, X., Wang, X., Shen, X., Chen, X., Zhang, X., Chen, X., Nie, X., Sun, X., Wang, X., Cheng, X., Liu, X., Xie, X., Liu, X., Yu, X., Song, X., Shan, X., Zhou, X., Yang, X., Li, X., Su, X., Lin, X., Li, Y. K., Wang, Y. Q., Wei, Y. X., Zhu, Y. X., Zhang, Y., Xu, Y., Xu, Y., Huang, Y., Li, Y., Zhao, Y., Sun, Y., Li, Y., Wang, Y., Yu, Y., Zheng, Y., Zhang, Y., Shi, Y., Xiong, Y., He, Y., Tang, Y., Piao, Y., Wang, Y., Tan, Y., Ma, Y., Liu, Y., Guo, Y., Wu, Y., Ou, Y., Zhu, Y., Wang, Y., Gong, Y., Zou, Y., He, Y., Zha, Y., Xiong, Y., Ma, Y., Yan, Y., Luo, Y., You, Y., Liu, Y., Zhou, Y., Wu, Z. F., Ren, Z. Z., Ren, Z., Sha, Z., Fu, Z., Xu, Z., Huang, Z., Zhang, Z., Xie, Z., Zhang, Z., Hao, Z., Gou, Z., Ma, Z., Yan, Z., Shao, Z., Xu, Z., Wu, Z., Zhang, Z., Li, Z., Gu, Z., Zhu, Z., Liu, Z., Li, Z., Xie, Z., Song, Z., Gao, Z., and Pan, Z. Deepseek-V3 Technical Report. 2024. doi: 10.48550/ARXIV.2412.19437. URL <https://arxiv.org/abs/2412.19437>.
- Dettmers, T., Pagnoni, A., Holtzman, A., and Zettlemoyer, L. Qlora: Efficient finetuning of quantized llms. *Advances in neural information processing systems*, 36:10088–10115, 2023.
- Ding, H., Li, Y., Wang, J., and Chen, H. Large Language Model Agent in Financial Trading: A Survey. 2024. doi: 10.48550/ARXIV.2408.06361. URL <https://arxiv.org/abs/2408.06361>.

- Dong, Y., Wu, F., Zhang, K., Dai, Y., Zhang, S., Ye, W., Chen, S., and Cheng, Z.-Q. Large Language Model Agents in Finance: A Survey Bridging Research, Practice, and Real-World Deployment. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pp. 17889–17907. Association for Computational Linguistics, 2025. doi: 10.18653/v1/2025.findings-emnlp.972. URL <http://dx.doi.org/10.18653/v1/2025.findings-emnlp.972>.
- Gunasekar, S., Zhang, Y., Aneja, J., Mendes, C. C. T., Del Giorno, A., Gopi, S., Javaheripi, M., Kauffmann, P., de Rosa, G., Saarikivi, O., et al. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*, 2023.
- Hao, S., Gu, Y., Ma, H., Hong, J., Wang, Z., Wang, D., and Hu, Z. Reasoning with language model is planning with world model. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 8154–8173, 2023.
- Hsieh, C.-Y., Li, C.-L., Yeh, C.-K., Nakhost, H., Fujii, Y., Ratner, A., Krishna, R., Lee, C.-Y., and Pfister, T. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 8003–8017, 2023.
- Huang, A. H., Wang, H., and Yang, Y. Finbert: A large language model for extracting information from financial text. *Contemporary Accounting Research*, 40(2):806–841, 2023.
- Hurst, A., Lerer, A., Goucher, A. P., Perelman, A., Ramesh, A., Clark, A., Ostrow, A., Welihinda, A., Hayes, A., Radford, A., et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. d. I., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mistral 7b. 2023. doi: 10.48550/ARXIV.2310.06825. URL <https://arxiv.org/abs/2310.06825>.
- Krumdick, M., Koncel-Kedziorski, R., Lai, V. D., Reddy, V., Lovering, C., and Tanner, C. Bizbench: A quantitative reasoning benchmark for business and finance. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8309–8332, 2024.
- Li, X., Li, Z., Shi, C., Xu, Y., Du, Q., Tan, M., and Huang, J. Alphafin: Benchmarking financial analysis with retrieval-augmented stock-chain framework. In *Proceedings of the 2024 joint international conference on computational linguistics, language resources and evaluation (LREC-COLING 2024)*, pp. 773–783, 2024.
- Li, Y., Wang, S., Ding, H., and Chen, H. Large Language Models in Finance: A Survey. In *4th ACM International Conference on AI in Finance*, pp. 374–382. ACM, nov 25 2023. doi: 10.1145/3604237.3626869. URL <http://dx.doi.org/10.1145/3604237.3626869>.
- Lightman, H., Kosaraju, V., Burda, Y., Edwards, H., Baker, B., Lee, T., Leike, J., Schulman, J., Sutskever, I., and Cobbe, K. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.
- Liu, Z., Huang, D., Huang, K., Li, Z., and Zhao, J. Finbert: A pre-trained financial language representation model for financial text mining. In *Proceedings of the twenty-ninth international conference on international joint conferences on artificial intelligence*, pp. 4513–4519, 2021.
- Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegrefe, S., Alon, U., Dziri, N., Prabhumoye, S., Yang, Y., et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36: 46534–46594, 2023.
- Meta. Meta llama 3.1 model information, Jul 2024. URL https://github.com/meta-llama/llama-models/blob/main/models/llama3_1/MODEL_CARD.md.
- Meta. Llama 4 scout model card, April 2025. URL https://github.com/meta-llama/llama-models/blob/main/models/llama4/MODEL_CARD.md.
- OpenAI et al., Nov 2025. URL <https://openai.com/index/gpt-5-1/>.

- Shinn, N., Cassano, F., Gopinath, A., Narasimhan, K., and Yao, S. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023.
- Turpin, M., Michael, J., Perez, E., and Bowman, S. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. *Advances in Neural Information Processing Systems*, 36:74952–74965, 2023.
- Uesato, J., Kushman, N., Kumar, R., Song, F., Siegel, N., Wang, L., Creswell, A., Irving, G., and Higgins, I. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- Wang, L., Du, Y., Lin, J., Chen, K., and Liew, S. C. Rephrase and contrast: Fine-tuning language models for enhanced understanding of communication and computer networks. In *2025 International Conference on Computing, Networking and Communications (ICNC)*, pp. 588–594, 2025. doi: 10.1109/ICNC64010.2025.10993716.
- Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., and Zhou, D. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- Webb, T., Holyoak, K. J., and Lu, H. Emergent analogical reasoning in large language models. *Nature Human Behaviour*, 7(9):1526–1541, 2023.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- 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.
- Xiao, Y., Sun, E., Luo, D., and Wang, W. Tradingagents: Multi-agents llm financial trading framework, 2025. URL <https://arxiv.org/abs/2412.20138>.
- Xie, Q., Han, W., Chen, Z., Xiang, R., Zhang, X., He, Y., Xiao, M., Li, D., Dai, Y., Feng, D., et al. Finben: A holistic financial benchmark for large language models. *Advances in Neural Information Processing Systems*, 37: 95716–95743, 2024.
- Yang, A., Li, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Gao, C., Huang, C., Lv, C., et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Yang, H., Liu, X.-Y., and Wang, C. D. Fingpt: Open-source financial large language models. *FinLLM Symposium at IJCAI 2023*, 2023a.
- Yang, Y., Tang, Y., and Tam, K. Y. Investlm: A large language model for investment using financial domain instruction tuning, 2023b.
- Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
- Yu, Y., Yao, Z., Li, H., Deng, Z., Jiang, Y., Cao, Y., Chen, Z., Suchow, J., Cui, Z., Liu, R., et al. Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making. *Advances in Neural Information Processing Systems*, 37:137010–137045, 2024.
- Zelikman, E., Wu, Y., Mu, J., and Goodman, N. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- Zelikman, E., Harik, G., Shao, Y., Jayasiri, V., Haber, N., and Goodman, N. D. Quiet-star: Language models can teach themselves to think before speaking. *arXiv preprint arXiv:2403.09629*, 2024.
- Zhang, B., Yang, H., and Liu, X.-Y. Instruct-fingpt: Financial sentiment analysis by instruction tuning of general-purpose large language models. *FinLLM Symposium at IJCAI 2023*, 2023.

Zhang, C., Liu, X., Jin, M., Zhang, Z., Li, L., Wang, Z., Hua, W., Shu, D., Zhu, S., Jin, X., et al. When ai meets finance (stockagent): Large language model-based stock trading in simulated real-world environments. *arXiv preprint arXiv:2407.18957*, 2024.

Zhao, Y., Li, Y., Li, C., and Zhang, R. MultihierTT: Numerical reasoning over multi hierarchical tabular and textual data. *arXiv preprint arXiv:2206.01347*, 2022.

Zhou, T., Wang, P., Wu, Y., and Yang, H. Finrobot: AI agent for equity research and valuation with large language models. In *ICAIF 2024: The 1st Workshop on Large Language Models and Generative AI for Finance*, 2024.