# Look-Ahead-Bench: a Standardized Benchmark of Look-ahead Bias in Point-in-Time LLMs for Finance

**Mostapha Benhenda** [*]

## Abstract

We introduce Look-Ahead-Bench, a standardized benchmark measuring look-ahead bias in Point-in-Time (PiT) Large Language Models (LLMs) within realistic and practical financial workflows. Unlike most existing approaches that primarily test inner lookahead knowledge via Q&A, our benchmark evaluates model behavior in practical scenarios. To distinguish genuine predictive capability from memorization-based performance, we analyze performance decay across temporally distinct market regimes, incorporating several quantitative baselines to establish performance thresholds.

We evaluate prominent open-source LLMs—Llama 3.1 (8B and 70B) and DeepSeek 3.2—against a family of Point-in-Time LLMs (Pitinf-Small, Pitinf-Medium, and frontier-level model Pitinf-Large) from PiT-Inference. Results reveal significant lookahead bias in standard LLMs, as measured with alpha decay, unlike Pitinf models, which demonstrate improved generalization and reasoning abilities as they scale in size. This work establishes a foundation for the standardized evaluation of temporal bias in financial LLMs and provides a practical framework for identifying models suitable for real-world deployment. Code is available on GitHub: https://github.com/benstaf/lookaheadbench

## 1. Introduction

Large Language Models are revolutionizing quantitative finance. FinGPT (Yang et al., 2023), BloombergGPT (Wu et al., 2023), FinMem (Yu et al., 2023), FinRL-DeepSeek (Benhenda, 2025), and TradingAgents (Xiao et al., 2025), along with multi-agent architectures like Hedge-Agents (Li et al., 2025a) and the AI Hedge Fund (Singh, 2025), have demonstrated remarkable capabilities in processing financial data, generating trading signals, and executing complex investment strategies. However, beneath these promising results lies a fundamental methodological challenge that threatens the validity of LLM-based financial applications: look-ahead bias.

Look-ahead bias occurs when models access information that would not have been available at the time of prediction, creating artificially inflated performance metrics that evaporate in real-world deployment. Consequently, back-tested returns often collapse once the model's knowledge window ends and trading gets into genuinely unknown territory (Li et al., 2025b).

### 1.1. The Challenge of Temporal Bias in Financial LLMs

LLMs are pre-trained on web-scale corpora containing extensive financial data, including post-hoc explanations of market events, historical price movements, and retrospective analyses. When an LLM encounters a prompt about "NVIDIA's performance in 2023," it may have been trained on text explicitly stating "NVIDIA surged 190% in 2023 on AI boom" (Li et al., 2025b). In this case, the model does not learn predictive relationships; it memorizes outcomes and recites them during evaluation.

This pre-training contamination creates systematic biases that are particularly pernicious in finance, where temporal causality is paramount. Several studies have documented various manifestations of this issue:

- **Training Data Leakage:** Models recall specific stock prices, earnings figures, and market events directly from their training corpus (Lopez-Lira et al., 2025b).

- **Entity Memorization:** Even with masked identifiers, models can recognize companies from contextual clues and retrieve associated historical performance (Sarkar & Vafa, 2024).

- **Temporal Reasoning Failures:** Models struggle to distinguish between information available at different time points, leading to logically impossible predictions (Yan & Tang, 2026; Gao et al., 2025).

. Correspondence to: Mostapha Benhenda <mostaphabenhenda@gmail.com>.

## 1.2. Limitations of Current Evaluation Approaches

Existing approaches to detecting lookahead bias have focused primarily on testing inner knowledge through question-answering tasks, or evaluating the memorization of specific facts. While valuable, these evaluation methods may not reflect the practical consequences of bias in real-world financial settings.

The pragmatic test of a financial model's capability does not lie in its ability (or lack thereof) to recall historical patterns from its training data, but in its practical capacity to generalize to genuinely novel market conditions.

## 1.3. Our Contributions

We introduce Look-Ahead-Bench, a standardized benchmark that addresses these limitations through three key contributions:

1. **Practical Trading Workflow Integration:** Unlike knowledge-based benchmarks, our evaluation uses realistic agentic trading systems that make actual portfolio decisions. We utilize the AI Hedge Fund framework (Singh, 2025), a widely adopted open-source project (+45k stars on GitHub), ensuring evaluation moves beyond toy problems to assess real-world deployment viability.

2. **Point-In-Time LLMs Evaluation:** We explicitly evaluate Pitinf models—commercial LLMs specifically designed to effectively remove lookahead bias— alongside standard foundation models.

3. **Alpha Decay Metric:** We introduce a rigorous metric for measuring bias by quantifying the performance drop between in-sample (training window) and out-of-sample (future) periods.

Our results demonstrate that standard LLMs exhibit significant lookahead bias, with alpha decay exceeding -15 percentage points between periods, while Pitinf models maintain stable performance. This work provides a diagnostic tool for identifying bias in financial LLMs.

## 2. Related Work

### 2.1. Lookahead Bias in Financial LLMs

The recognition of lookahead bias as a critical issue in financial LLM applications has grown rapidly in recent years. Glasserman & Lin (2023) provided one of the earliest systematic assessments, investigating both lookahead bias and distraction effects in GPT-based sentiment analysis for stock return prediction. Their work revealed that anonymizing company identifiers could actually improve performance by reducing distraction effects, while highlighting the challenge of disentangling genuine sentiment analysis from memorized associations.

Levy (2025) emphasized implications for financial applications and proposed mitigation strategies including careful data cutoff management and entity embedding neutralization. Gao et al. (2025) developed a statistical test for detecting lookahead bias using Lookahead Propensity (LAP), a measure based on membership inference attack techniques. Their approach correlates model familiarity with training data (estimated through token probability analysis) with forecast accuracy. While innovative, this method focuses on individual prediction tasks rather than end-to-end trading system evaluation.

Other contributions have expanded both the empirical and methodological study of lookahead bias in financial language models, proposing standardized diagnostics, benchmark-style evaluations, and controlled temporal stress tests that reveal performance inflation caused by future data contamination in forecasting and trading settings (Noguer i Alonso, 2024; Lopez-Lira et al., 2025a).

### 2.2. Information Leakage and Memorization

The broader issue of information leakage in LLMs has received extensive attention. Li et al. (2025b) introduced FinLeak-Bench and the FactFin framework to systematically quantify information leakage across four dimensions, demonstrating that most published LLM-based agents fail to beat baselines once their knowledge cutoff is passed.

Lopez-Lira et al. (2025b) documented the "memorization problem" through extensive testing of GPT-4o's recall capabilities, showing that models can recall exact S&P 500 closing prices with less than 1% error for dates within their training window, while errors "explode" for post-cutoff dates. Carlini et al. (2021; 2022) established theoretical and empirical foundations for understanding memorization in large language models, demonstrating that training data can be extracted from models through carefully designed queries.

Complementary studies further show that pretraining contamination and implicit temporal leakage can persist even without explicit date cues, and that leakage effects scale with model size unless explicitly mitigated (Kong et al., 2025; Drinkall et al., 2024).

### 2.3. Point-in-Time and Temporally Aware Models

The development of Point-in-Time (PiT) models represents a direct response to lookahead bias concerns. Models like Time Machine GPT (Drinkall et al., 2024), ChronoGPT (He et al., 2025) and DatedGPT (Yan & Tang, 2026) introduce time-aware frameworks that train models strictly on

pre-cutoff data to ensure temporal integrity. Their approach demonstrates that explicit temporal constraints can prevent future information leakage. Merchant & Levy (2025) proposed logit adjustment from specialized auxiliary models, allowing both verbatim and semantic knowledge to be removed without retraining large base models.

There are also proprietary Point-In-Time LLMs designed for backtesting, such as those from PiT-Inference (PiT-Inference, 2026). These are available in three sizes: Small size ($<$10B parameters) for low-latency, Medium size ($<$100B parameters), and Large size ($>$500B parameters) for frontier-grade reasoning tasks.

## 2.4. Agentic Trading Systems

The application of LLMs to agentic trading has produced numerous innovative systems. Yu et al. (2023; 2024) developed FinMem and FinCon, demonstrating sophisticated memory and multi-agent capabilities. Zhang et al. (2024b) created a multimodal foundation agent with tool augmentation, while Yang et al. (2024) introduced FinRobot for equity research and valuation.

The AI Hedge Fund open-source project (Singh, 2025) provides LLM-based trading agents and serves as the foundation for our benchmark implementation, due to its significant impact (+45k stars on GitHub) on the community.

## 3. Methodology

### 3.1. Dual-Period Design

Our evaluation compares model performance across two carefully selected periods to test models with knowledge cutoffs through 2023-2024, while maintaining similar market characteristics for fair comparison, like in FinLeak-Bench (Li et al., 2025b):

**Period P1 (In-Sample):** April 1, 2021 – September 30, 2021

- Duration: 6 months
- Buy-and-Hold Return: +25.32%
- Purpose: Establish baseline performance within the potential training window of Llama 3.1 (8B, 70B) and DeepSeek 3.2.

**Period P2 (Out-of-Sample):** July 1, 2024 – December 31, 2024

- Duration: 6 months
- Market Regime: AI-driven rally with mixed sentiment.

- Buy-and-Hold Return: +24.75% (similar to P1).
- Purpose: Test generalization to post-cutoff periods for Llama 3.1 (8B, 70B) and DeepSeek 3.2.

### 3.2. Portfolio and Trading Universe

Our benchmark employs a focused trading universe consisting of five large-cap technology stocks: AAPL (Apple Inc.), MSFT (Microsoft Corporation), GOOGL (Alphabet Inc.), NVDA (NVIDIA Corporation), and TSLA (Tesla Inc.).

### 3.3. Evaluation Metrics

#### 3.3.1. ALPHA CALCULATION

For each model and period, we calculate alpha as:

$$\alpha = R_{\text{LLM}} - R_{\text{Buy\&Hold}} \tag{1}$$

This measures the excess return generated by the LLM strategy relative to a passive buy-and-hold benchmark.

#### 3.3.2. ALPHA DECAY ANALYSIS

The core bias indicator is alpha decay between periods:

$$\alpha_{\text{Decay}} = \alpha_{P2} - \alpha_{P1} \tag{2}$$

Negative alpha decay indicates that the model's relative performance deteriorated in the out-of-sample period, suggesting potential lookahead bias in the in-sample results.

### 3.4. Quantitative Baselines

We benchmark the AI agents against a suite of traditional quantitative strategies. These strategies span passive, systematic, trend-following, and contrarian approaches, providing a comprehensive reference for evaluating relative performance and robustness across market regimes. All strategies are executed with the same initial capital, fractional shares allowed, and monthly rebalancing (except where noted).

#### 3.4.1. BUY & HOLD (PASSIVE)

The simplest baseline allocates capital equally across all tickers at the start of the period and holds the positions without any rebalancing. This represents a passive, long-only exposure to the selected universe and serves as the reference for computing alpha (excess return over this benchmark).

#### 3.4.2. EQUAL-WEIGHT MONTHLY REBALANCE (SYSTEMATIC)

This strategy maintains equal dollar allocation to each ticker by rebalancing at the start of each month. It captures

broad market exposure while controlling concentration risk, often outperforming naive buy-and-hold in diversified universes due to the rebalancing premium.

### 3.4.3. MOMENTUM (3-MONTH, TOP HALF)

A classic trend-following strategy that, at each monthly rebalance, ranks tickers by their past 3-month total return and allocates equally to the top half (or all tickers equally if insufficient history). This exploits persistence in intermediate-term price trends and is one of the most robust quantitative factors historically.

### 3.4.4. MEAN REVERSION (3-MONTH, BOTTOM HALF)

The contrarian counterpart to momentum: at each monthly rebalance, the strategy ranks tickers by past 3-month return and allocates equally to the bottom half of performers (or all equally if insufficient history). It bets on price reversals and tends to perform well in range-bound or high-dispersion regimes but can suffer in strong trends.

### 3.4.5. MOVING AVERAGE CROSSOVER (50/100-DAY, DAILY REBALANCE)

A trend-following rule that holds a ticker only when its 50-day simple moving average is above its 100-day simple moving average. Positions are rebalanced daily: capital is equally allocated to all tickers meeting the condition (cash otherwise). This classic dual-moving-average system aims to ride trends while avoiding major drawdowns.

### 3.4.6. RANDOM NOISE (CONTROL)

A deliberately naive control strategy that, each day, applies random weights (±40% long, ±40% short, 20% zero per ticker, normalized) to daily ticker returns. This "monkey with a dartboard" baseline helps gauge whether observed performance exceeds what pure randomness could achieve in the same universe.

### 3.5. LLM-based Trading Agents

We implement the LLM-based trading agents using the open-source AI Hedge Fund framework (Singh, 2025), a popular proof-of-concept system (45.3k stars on GitHub as of January 2026) for exploring AI-driven investment decisions. The repository is available at `https://github.com/virattt/ai-hedge-fund`.

The framework employs a multi-agent architecture where specialized LLM-powered agents—including a Valuation Agent, Sentiment Agent, Fundamentals Agent, Technicals Agent, and others modeled after prominent investors (e.g., Warren Buffett)—analyze market data and generate trading signals. A Risk Manager agent computes position limits based on risk metrics, and a Portfolio Manager agent syn-

thesizes the signals into final allocation decisions. Trades are simulated only (no live execution), making it ideal for research and benchmarking.

Key features relevant to our evaluation include:

- Support for multiple LLM providers via API (OpenAI, Groq, Anthropic, DeepSeek)

- Built-in access to historical price data for popular tickers (AAPL, GOOGL, MSFT, NVDA, TSLA) without requiring an external API key—aligning precisely with our trading universe.

In our experiments, we adapt the framework to a consistent monthly rebalancing protocol (aligned with the quantitative baselines): at the start of each month, the agents are prompted with historical data up to that point, generate target weights for the five stocks, and the weights are normalized and executed (fractional shares allowed, equal initial capital).

We evaluate the following models within this agentic workflow:

**Standard Foundation Models**   Widely used open-source LLMs with potential training data contamination extending into 2023–2024:

- Meta Llama 3.1 (8B and 70B parameters), with cutoff date of December 2023

- DeepSeek 3.2, with cutoff date of July 2024

**Point-in-Time (PiT) Models**   The proprietary Pitinf family from PiT-Inference (PiT-Inference, 2026), designed with effective temporal cutoffs of January 2020 to eliminate lookahead bias:

- Pitinf-Small (<10B parameters) — suited for low-latency tasks.

- Pitinf-Medium (<100B parameters).

- Pitinf-Large (>500B parameters) — for frontier-level reasoning performance.

## 4. Experimental Results

Table 1 summarizes the performance of quantitative baselines and AI models across both periods.

### 4.1. Key Insights

**Relative to Quant Strategies:**

*Table 1.* Performance Comparison across quantitative baselines and AI models. Note the "Scaling Paradox" in P2: while Standard models degrade as they scale (due to stronger false priors), PiT models improve as they scale (due to cleaner reasoning).

| Model / Strategy | Variant | P1 Return (%) | P1 Alpha (pp) | P2 Return (%) | P2 Alpha (pp) | Alpha Decay |
|---|---|---|---|---|---|---|
| *QUANT STRATEGIES (Baseline)* | | | | | | |
| Buy & Hold | Passive | +25.32 | +0.00 | +24.75 | +0.00 | +0.00 |
| Equal Weight | Systematic | +25.68 | +0.36 | +22.33 | -2.42 | -2.78 |
| Momentum (3M) | Systematic | +33.28 | +7.96 | +30.50 | +5.75 | -2.21 |
| Mean Reversion | Systematic | +20.70 | -4.62 | +9.35 | -15.40 | -10.78 |
| MA Crossover | Trend | -2.46 | -27.78 | +6.91 | -17.84 | +9.94 |
| Random Noise | Control | +11.43 | -13.89 | +0.56 | -24.19 | -10.30 |
| *AI MODELS (Agents)* | | | | | | |
| Llama 3.1 8B | Standard | +39.13 | +13.81 | +21.33 | -3.42 | -17.23 |
| Llama 3.1 70B | Standard | +44.59 | +19.27 | +28.77 | +4.02 | -15.25 |
| DeepSeek 3.2 | Standard | +46.05 | +20.73 | +23.71 | -1.04 | -21.77 |
| Pitinf-Small | PiT | +25.07 | -0.25 | +24.81 | +0.06 | +0.31 |
| Pitinf-Medium | PiT | +27.76 | +2.44 | +28.04 | +3.29 | +0.85 |
| Pitinf-Large | PiT | +31.34 | +6.02 | +32.07 | +7.32 | +1.30 |

- **Momentum Robustness:** The best quantitative strategy (Momentum) generated +7.96pp alpha in P1 and maintained +5.75pp in P2, showing only mild decay (-2.21pp).

- **Mean Reversion Failure:** The worst strategy (Mean Reversion) decayed significantly (-10.78pp), highlighting the risk of overfitting strategies to specific mean-reverting regimes in trending markets.

**The Scaling Paradox and Inverse Scaling:** Our results highlight a critical divergence in how model scaling affects performance in the presence of lookahead bias. We observe what we term the "Scaling Paradox":

- **Standard Models (The "Memory Trap"):** For standard models, scaling size actually *hurt* P2 performance relative to the peak. While DeepSeek 3.2 achieved the highest P1 Alpha (+20.73pp) due to superior memorization of 2021 data, it suffered the most severe collapse in P2 (-21.77pp decay).

- *Interpretation:* This aligns with the "Inverse Scaling" phenomenon (McKenzie et al., 2023), where performance on specific tasks degrades as model scale increases. In financial forecasting, larger models develop stronger, more brittle priors based on training data. As noted by Lopez-Lira et al. (2025c), larger models are more adept at recalling exact historical values (perfect memory), but this "photographic" recall becomes a liability during regime shifts. When deployed in P2 (2024), these strong internal priors conflict with reality, leading to confident hallucinations that override immediate contextual signals.

- **Point-in-Time Models (The "Reasoning Dividend"):** Conversely, the Pitinf family demonstrates a positive scaling law in the absence of bias. Once the distractor of "future memory" is removed, scaling the model size improves genuine financial reasoning—such as sentiment interpretation and macro-analysis—rather than merely scaling the capacity to memorize historical price paths.

## 5. Conclusion and Future Work

This paper is a proof-of-concept of how Look-Ahead-Bench can serve as a diagnostic tool for look-ahead bias in financial LLMs. Standard foundation models, while powerful, suffer from severe alpha decay when moved from memorized training windows to new data. In contrast, Pitinf models demonstrate better generalization.

However, more experimentations would be needed, such as:

### 5.1. Expand Experimental Scope

The current evaluation is restricted to five large-cap US technology stocks. To ensure robustness, more diverse backtesting settings are needed, with at least 20–30 stocks across multiple sectors, including:

- Small-cap stocks: These are less likely to be memorized in training corpora.

- Non-tech sectors: E.g., Healthcare, Industrials.

- Diverse Assets: Commodities or FX, though LLM applicability there differs.

- Multiple Time Periods: Testing 4+ periods of varying lengths.

## 5.2. Expand Models and Agents

Look-ahead Bench needs to be enriched with a broader range of trading agents, such as FinMem (Yu et al., 2023), FinGPT (Yang et al., 2023), FinRL-DeepSeek (Benhenda, 2025), TradingAgents (Xiao et al., 2025), and Hedge-Agents (Li et al., 2025a).

Additionally, the scope of autonomous trading architectures should be expanded by integrating specialized agents including FinAgent (Zhang et al., 2024b), FinRobot (Yang et al., 2024), and StockAgent (Zhang et al., 2024a). Furthermore, QuantAgent (Xiong et al., 2025) needs to be incorporated for quantitative investment strategies, while CryptoTrade (Li et al., 2024) should be added to extend coverage to cryptocurrency markets.

To assess fundamental analysis and predictive capabilities, models focusing on factor discovery and market forecasting should be considered. This should include LLMFactor (Cao, 2025) for quantitative factor mining, as well as predictive agents such as TradingGPT (Li et al., 2023), AlphaGPT (Wang et al., 2025), MarketGPT (Wheeler & Varner, 2024), and StockGPT (Mai, 2024). Evaluation frameworks like PIXIU (Xie et al., 2023) and FinPos (Bijia Liu and Ronghao Dang, 2025) should also be utilized for their comprehensive datasets and insights into position sizing.

Finally, to address complex market dynamics involving interaction and high-frequency trading, the benchmark needs to be extended to support multi-agent reinforcement learning environments. This should involve frameworks such as Language Model Guided RL (Darmanin & Vella, 2025). Recent tools like AutoTrader (RexiaAI, 2025) should also be examined to cover the spectrum from research prototypes to deployable trading systems.

## 5.3. Expand Backtesting Methodologies

The benchmark should aim to move beyond historical price paths by employing advanced validation techniques, including:

- Synthetic Data and Counterfactuals: Testing how models react to market events that could have happened but didn't.

- Rademacher Anti-Serum (RAS): Utilizing methods like those proposed by Paleologo (2025) to rigorously test for backtesting overfitting.

## References

Benhenda, M. FinRL-DeepSeek: LLM-infused risk-sensitive reinforcement learning for trading agents, 2025. URL https://arxiv.org/abs/2502.07393.

Bijia Liu and Ronghao Dang. FinPos: A position-aware trading agent system for real financial markets, 2025. URL https://arxiv.org/abs/2510.27251.

Cao, L. Chain-of-alpha: Unleashing the power of large language models for alpha mining in quantitative trading, 2025. URL https://arxiv.org/abs/2508.06312.

Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, U., Oprea, A., and Raffel, C. Extracting training data from large language models. In *USENIX Security Symposium*, 2021.

Carlini, N., Ippolito, D., Jagielski, M., Lee, K., Tramer, F., and Zhang, C. Quantifying memorization across neural language models. In *International Conference on Machine Learning (ICML)*, 2022.

Darmanin, A. and Vella, V. Language model guided reinforcement learning in quantitative trading, 2025. URL https://arxiv.org/abs/2508.02366.

Drinkall, F., Rahimikia, E., Pierrehumbert, J. B., and Zohren, S. Time Machine GPT, 2024. URL https://arxiv.org/abs/2404.18543.

Gao, Z., Jiang, W., and Yan, Y. A test of lookahead bias in LLM forecasts, 2025. URL https://arxiv.org/abs/2512.23847.

Glasserman, P. and Lin, C. Assessing Look-Ahead Bias in Stock Return Predictions Generated By GPT Sentiment Analysis, 2023. URL https://arxiv.org/abs/2309.17322.

He, S., Lv, L., Manela, A., and Wu, J. Chronologically consistent large language models, 2025. URL https://arxiv.org/abs/2502.21206.

Kong, Y., Hwang, Y., Kaiser, M., Vryonides, C., Oomen, R., and Zohren, S. Fusing narrative semantics for financial volatility forecasting, 2025. URL https://arxiv.org/abs/2510.20699.

Levy, B. Caution Ahead: Numerical Reasoning and Look-ahead Bias in AI Models. *SSRN Electronic Journal*, 2025. URL https://ssrn.com/abstract=5082861.

Li, X., Zeng, Y., Xing, X., Xu, J., and Xu, X. HedgeAgents: A Balanced-aware Multi-agent Financial Trading System, 2025a. URL https://arxiv.org/abs/2502.13165.

Li, X., Zeng, Y., Xing, X., Xu, J., and Xu, X. Profit Mirage: Revisiting Information Leakage in LLM-based Financial Agents, 2025b. URL https://arxiv.org/abs/2510.07920.

Li, Y., Yu, Y., Li, H., Chen, Z., and Khashanah, K. TradingGPT: Multi-Agent System with Layered Memory and Distinct Characters for Enhanced Financial Trading Performance, 2023. URL https://arxiv.org/abs/2309.03736.

Li, Y., Luo, B., Wang, Q., Chen, N., Liu, X., and He, B. A Reflective LLM-based Agent to Guide Zero-shot Cryptocurrency Trading, 2024. URL https://arxiv.org/abs/2407.09546.

Lopez-Lira, A., Choi, C., Kim, Y., Kwon, J., Kim, J., and Yun, S. Thematic Scoring: Quantifying Contextual Narratives using Language Models. *SSRN Electronic Journal*, 2025a. URL https://ssrn.com/abstract=5233994. SSRN 5233994.

Lopez-Lira, A., Tang, Y., and Zhu, M. The Memorization Problem: Can We Trust LLMs' Economic Forecasts? *SSRN Electronic Journal*, 2025b. URL https://ssrn.com/abstract=5217505.

Lopez-Lira, A., Tang, Y., and Zhu, M. The Memorization Problem: Can We Trust LLMs' Economic Forecasts? *arXiv preprint arXiv:2504.14765*, 2025c. URL https://arxiv.org/abs/2504.14765.

Mai, D. StockGPT: A GenAI Model for Stock Prediction and Trading, 2024. URL https://arxiv.org/abs/2404.05101.

McKenzie, I. R., Lyzhov, A., Pieler, M., Parrish, A., Mueller, A., Prabhu, A., McLean, E., Kirtland, A., Ross, A., Liu, A., et al. Inverse Scaling: When Bigger Isn't Better. *Transactions on Machine Learning Research*, 2023. URL https://openreview.net/forum?id=DwgRm72GQF.

Merchant, H. and Levy, B. A fast and effective solution to the problem of look-ahead bias in LLMs, 2025. URL https://arxiv.org/abs/2512.06607.

Noguer i Alonso, M. Look-Ahead Bias in Large Language Models (LLMs): Implications and Applications in Finance. *SSRN Electronic Journal*, 2024. URL https://ssrn.com/abstract=5022165. SSRN 5022165.

Paleologo, G. A. *The Elements of Quantitative Investing*. John Wiley & Sons, 2025.

PiT-Inference. Introducing Pitinf models, 2026.

RexiaAI. AutoTrader: An automated trading framework using LLMs. https://github.com/RexiaAI/AutoTrader, 2025.

Sarkar, A. and Vafa, K. Lookahead Bias in Pretrained Language Models, 2024. URL https://ssrn.com/abstract=4754678. SSRN 4754678.

Singh, V. AI hedge fund: Open-source framework for LLM-based trading agents. https://github.com/virattt/ai-hedge-fund, 2025.

Wang, S., Yuan, H., Zhou, L., Ni, L. M., Shum, H.-Y., and Guo, J. Alpha-GPT: Human-AI Interactive Alpha Mining for Quantitative Investment, 2025. URL https://arxiv.org/abs/2308.00016.

Wheeler, A. and Varner, J. MarketGPT: Developing a Pre-trained transformer (GPT) for Modeling Financial Time Series, 2024. URL https://arxiv.org/abs/2411.16585.

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, 2023. URL https://arxiv.org/abs/2303.17564.

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., Zhang, X., Lai, Y., Peng, M., Lopez-Lira, A., and Huang, J. PIXIU: A Large Language Model, Instruction Data and Evaluation Benchmark for Finance, 2023. URL https://arxiv.org/abs/2306.05443.

Xiong, F., Zhang, X., Feng, A., Sun, S., and You, C. QuantAgent: Price-Driven Multi-Agent LLMs for High-Frequency Trading, 2025. URL https://arxiv.org/abs/2509.09995.

Yan, Y. and Tang, R. DatedGPT: Preventing Lookahead Bias in Large Language Models with Time-Aware Pretraining. *AFA Poster Session (2026)*, 2026.

Yang, H., Liu, X.-Y., and Wang, C. D. FinGPT: Open-source financial large language models, 2023. URL https://arxiv.org/abs/2306.06031.

Yang, H., Zhang, B., Wang, N., Guo, C., Zhang, X., Lin, L., Wang, J., Zhou, T., Guan, M., Zhang, R., and Wang, C. D. Finrobot: An open-source AI agent platform for financial applications using large language models, 05 2024. URL https://arxiv.org/abs/2405.14767.

Yu, Y., Li, H., Chen, Z., Jiang, Y., Li, Y., Zhang, D., and Khashanah, K. FinMem: A performance-enhanced LLM trading agent with layered memory and character design, 2023. URL https://arxiv.org/abs/2311.13743.

Yu, Y., Yao, Z., Li, H., Deng, Z., Cao, Y., Chen, Z., and Xie, Q. FinCon: A synthesized LLM multi-agent system with conceptual verbal reinforcement for enhanced financial decision making, 2024. URL https://arxiv.org/abs/2407.06567.

Zhang, C., Liu, X., Zhang, Z., Jin, M., Li, L., Wang, Z., Hua, W., Shu, D., Zhu, S., Jin, X., Li, S., Du, M., and Zhang, Y. When AI meets finance (Stock-Agent): Large language model-based stock trading in simulated real-world environments, 07 2024a. URL https://arxiv.org/abs/2407.18957.

Zhang, W., Zhao, L., Xia, H., Sun, S., Sun, J., Qin, M., Li, X., Zhao, Y., Zhao, Y., Cai, X., Zheng, L., Wang, X., and An, B. A multi-modal foundation agent for financial trading: Tool-augmented, diversified, and generalist, 02 2024b. URL https://arxiv.org/abs/2402.18485.