# To Trade or Not to Trade:
# An Agentic Approach to Estimating Market Risk Improves Trading Decisions

Dimitrios Emmanoulopoulos[1,*]        Ollie Olby[2]        Justin Lyon[2]

Namid R. Stillman[2,*]

July 14, 2025

## Abstract

Large language models (LLMs) are increasingly deployed in agentic frameworks, in which prompts trigger complex tool-based analysis in pursuit of a goal. While these frameworks have shown promise across multiple domains including in finance, they typically lack a principled model-building step, relying instead on sentiment- or trend-based analysis. We address this gap by developing an agentic system that uses LLMs to iteratively discover stochastic differential equations for financial time series. These models generate risk metrics which inform daily trading decisions. We evaluate our system in both traditional backtests and using a market simulator, which introduces synthetic but causally plausible price paths and news events. We find that model-informed trading strategies outperform standard LLM-based agents, improving Sharpe ratios across multiple equities. Our results show that combining LLMs with agentic model discovery enhances market risk estimation and enables more profitable trading decisions.

**Keywords:** LLMs, agentic finance, AI agents, model discovery, market risk, stochastic differential equations, multi-agent system, Monte Carlo simulations, agentic trading

## 1 Introduction

$\mathscr{T}$ HE recent development of large language models (LLMs) capable of wider ranging technical tasks has significantly changed the scope of quantitative analysis that can be, at least partially, automated. Applications range from technical analysis of a company's fundamental value, wider market sentiment, factor analysis and most tasks involving some form of natural language processing (NLP) [1, 2]. The implications to trading systems will likely be a dramatic increase in the rate and volume of market insights that can be generated to inform decisions.

The overall capabilities of LLMs have dramatically increased over the last five years [3]. This has led to an increase in the number of LLMs available, both as proprietary models from frontier labs or as smaller models with open-weights which can be run locally. Given this, the influence of LLMs on trading decisions is expected to be varied and highly model specific. Early work is starting to compare and benchmark these models in tasks specific to financial applications, such as trading decisions, portfolio optimisation, and market analysis [4–10]. As the number of models increases, and their underlying strengths and weaknesses become more apparent, it is expected that different classes of pre-trained models will be more regularly deployed to achieve certain objectives [11, 12].

While these objectives are likely to be significantly linked to NLP-based tasks, such as text summarisation, analysis, and generation, recent LLM architectures give early evidence that more complex tasks can also be automated. These LLMs, such as the 'o' series from OpenAI or 'R1' from DeepSeek, generate 'reasoning' tokens which result in the model performing more in-context analysis of the generated output and has lead to improved performance over a number of key evaluation measures [13, 14]. Though it is widely considered that these LLMs are still significantly weaker than human reasoning

* Corresponding authors: dimitrios.emmanoulopoulos@barclays.com, namid@simudyne.com
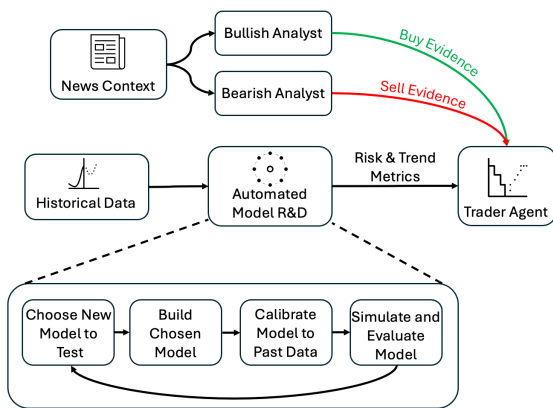
[1]Barclays, [2]Simudyne

Figure 1: *Agentic method for estimating market risk.* We identify a stochastic model of historical price paths using an agentic approach. Here, we follow a *builder-critic* pattern for model discovery, in which AI agents implement, simulate, test, and score proposed stochastic models. These models are generated as the output of our agents. We use this to construct and calibrate stochastic models of price paths which provide model-informed estimates for market risk. These risk metrics are provided to a *trader* agent which, along with market context, inform trading decisions. We find that our *risk analyst* agent, using this agentic model discovery loop, improves trading decisions.

at most tasks, they appear to demonstrate sufficient capabilities during in-context chain-of-thought messages that they are collectively referred to as 'reasoning' LLMs though the degree to which these LLMs are performing actual reasoning is a topic of active debate[3, 15, 16].

While the application of LLMs to finance tasks is still in its early stages, there is an open question around the extent in which LLMs will impact future trading decisions and whether complex tasks, such as model development, can even be automated. In this research study, we address this question. Specifically, we investigate the performance of LLMs when tasked with the identification of a realistic model of financial time series data and the subsequent impact of their analysis on trading decisions. In this sense, we make a separation of two key components within our agentic framework; the *analyst* and *trader* roles[1]. The *analyst* agents are capable of either calculating risk metrics (referred to as *risk analysts*) or summarising recent news stories for sentiment assessment (*news analysts*). Additionally, and perhaps most importantly, we focus on the development of the *risk analyst* to have sufficient capabilities to attempt model discovery. We structure this process of model discovery as the generation of a stochastic differential equation (SDE) that reproduces the statistics of a historical price path.

Our results show that LLMs are able to perform financial model discovery and refinement, through a process of implementing, calibrating, comparing, summarising, and scoring the proposed stochastic equations. We compare a subset of the latest LLMs on their ability to perform these tasks, highlighting those that are able to identify the stochastic models that best match historical data through a combination of statistical features, symbolic similarity and internal scoring. This agentic framework for model discovery is summarised as the *risk analyst* agent within our *trader-analyst* system.

Given the ability to identify the correct model formulation for the stochastic time series, we next investigate how these agents impact trading decisions, as measured by the impact of their quantification of market risk. To do so, we task a *risk analyst* with discovering suitable stochastic model descriptions for historical equity prices and then estimate market risk statistics (which we refer to as metrics) from these models. These risk metrics are then passed to a trader who can decide to buy or sell depending on these statistics as well as recent information about the stock, as outlined in Figure 1. We find that our model-informed *risk analyst* improves trading decisions and the degree of improvement is directly connected to the strength of the underlying LLM at model discovery. This highlights that model-informed risk analysis is expected to improve as LLM performance improves.

Finally, we expect that AI agents will be increasingly dependent on the underlying bias within their training data which consists of historical information up to the so-called knowledge cut-off data (the

---

[1]Throughout our manuscript, we denote LLM agents such as *builder*, *critic*, *analyast* and *trader*, by italics.

final date for which the training data relates to). This means that LLMs are likely to be ill-suited for out-of-distribution events that contradict their training data. To test this, we use the *intelligent backtesting* of the `Simudyne Horizon` market simulator[2]. This generates historical time series which matches the statistical properties of the price path but individual per-day shocks relate to synthetic macro-economic events. These results demonstrate that our findings are not dependent on training data, and are robust even if tested on 'causally consistent' synthetic data.

Overall, our results demonstrate that currently available LLMs are already sufficiently advanced to augment existing trading decisions with model-informed insights. Furthermore, the process of model discovery, previously assumed to require a high-degree of technical and market-specific knowledge, can be semi-automated in a principled way but we expect that this will continue to require some degree of human supervision. Overall, we find that these insights improve current trading decisions when tested using both traditional and intelligent backtesting methods. Finally, our research introduces new methods for evaluating the performance of LLMs. This form of evaluation highlights a novel new approach to benchmarking LLMs, where underlying LLMs are exposed to real-world tasks as in other evaluation frameworks focused on agentic systems [17].

## 1.1 Related work

$\mathscr{R}$ ECENT work has shown that LLMs display broad capabilities in many natural language tasks. These skills have been leveraged towards more complicated tasks through the use of agentic workflows. For example, recent work shows how to employ agentic LLM structures to develop trading strategies [18]. It uses specialist agents specifically tuned for tasks such as equity research, risk assessment, and strategy formulation to recreate the dynamics of a trading firm. This framework ensures precision and explainability in making financial trading decisions. Alternatively, a novel human-AI interface was developed for alpha mining [19]. By leveraging LLMs, 'Alpha-GPT' interprets quantitative researchers' ideas, generates formulaic alphas, and iteratively refines them through user feedback. By integrating natural language understanding, reasoning, and structured data synthesis, Alpha-GPT streamlines the discovery, evaluation, and deployment of novel trading signals.

Additionally, other work uses multi-modal LLMs to filter through eleven documents providing background in financial alpha mining strategies [20]. This filtered information is then passed to a multi-agent decision making framework to determine seed alphas. The performance from each agent is weighted using deep learning methods to determine the overall strategy. This is approach was expanded to consider a novel framework that leverages a mixture of experts (MoE) approach for stock trading using LLMs [21]. It integrates four LLM 'experts', each focusing on separate data sources: news articles, market data, alpha factors, and fundamental data. A general expert combines insights from the four experts for final predictions.

Recently, a layered memory approach was shown to enhance trading performance [22]. This leveraged LLMs' ability to determine the relevance of a particular piece of information and then store that in long, short or medium term memory. This memory storage framework was proposed to create a more human-like approach to trading, reflecting how individuals retain significant financial events over extended periods. Other work explores a three-step methodology that used LLMs to transform qualitative financial insights into quantitative predictions [23]. The LLM used a moving 5-shot prompting technique to reflect weekly trends and enhance the predictive abilities.

Finally, recent works explored model discovery using LLMs [24, 25]. These works used a builder-critic exploration loop, as first set-out in the seminal work by George E. Box [26]. This has been extended to financial markets by using vision language models which generate a neural-symbolic model for the fundamental price of an asset [25]. Our work builds upon this research by extending their neural-symbolic model discovery framework with LLMs to complement an agentic trading system.

---

[2]SIMUDYNE is a registered trademark of Simudyne Holdings Limited.

# 2 Methodology

$\mathscr{A}$ GENTIC AI workflows describe LLMs that use tools to achieve a designated task. In this research, we show how AI agents can be integrated into a model discovery loop to identify stochastic models for simulating financial time series. The resulting insights from the (LLM-discovered) stochastic simulator can inform trading decisions such as by determining model-specific risk metrics.

All of our experiments were performed on NVIDIA GPUs (V100s, A10s, A100s), for assessing LLM performance at discovering stochastic models. In total, we used ∼1,100 GPU hours to calibrate and simulate the proposed stochastic models. All calls to LLMs were made using `AWS Bedrock`, except for OpenAI which was made using the `OpenAI API`. For our experiments, we selected a subset of LLMs that were available during the period of January to May 2025. We acknowledge that our selection of language models is neither extensive nor exhaustive, and future work will extend our results to other large and multimodal language models. For data preprocessing, we used NVIDIA `RAPIDS`[3], for simulation and calibration of the stochastic models we used `diffrax` and `optax`, respectively [27], on NVIDIA GPUs. For LLM agentic prompt construction, we used `dspy` [28]. Finally, for benchmarking LLMs within an advanced backtesting environment, we used `Simudyne Horizon`, as outlined in Section 3.3.

## 2.1 Risk Analyst Agents

We combine an agentic AI framework with model discovery to identify stochastic models of financial time series. These models are used to simulate financial time series and then estimate model-informed risk statistics, such as value-at-risk (VaR) and maximum drawdown (MDD). These risk metrics are then used to inform an LLM-based trading agent.

### 2.1.1 Stochastic Simulator

We are specifically interested in identifying a model, $m \in \mathcal{M}$, where $m$ is a concrete realisation of the superset of stochastic processes, such as geometric Brownian motion. We constrain the space of all possible models, $\mathcal{M}$, to those that can be written as a stochastic differential equation (SDE), namely

$$\mathrm{d}S_t(\theta) = f(S_t, t; \theta)\,\mathrm{d}t + g(S_t, t; \theta) \circ \mathrm{d}W_t \tag{1}$$

in which $f(S_t, t)$ denotes the deterministic term, $g(S_t, t)$ denotes the diffusion term and we assume that $S_t$ is a continuous $\mathbb{R}^n$-valued process representing the evolution of price dynamics of an asset. The stochastic term includes composition with Brownian motion, $\circ\,\mathrm{d}W_t$, which we assume to be sampled from the normal distribution with mean zero and standard deviation proportional to the time step, $\mathrm{d}W_t \sim \mathcal{N}(0, \mathrm{d}t)$. This model, $m$, is then used to simulate the price dynamics of an asset, $S_t$, given some historical price data, $S_t^*$, and model-specific parameters, $\theta$.

The functional form of an SDE approximating financial time series is likely to be complex and time varying. This is because exogenous shocks to markets, such as the COVID-19 pandemic, are unlikely to be well-described when using historical data that does not contain these events. This is the 'black swan' problem of model selection in finance [29]. We note that functions $f$ and $g$ can be multi-dimensional complex functions but, in this work, we restrict the space of SDEs to a single independent variable, considering only the evolution of the price, $S_t$, and not, for example stochastic models of both volatility and price paths.

Given a proposed model of the market dynamics, $m$, we need to identify the correct parameters, $\theta$ that match model output, $S_t$, to historical data, $S_t^*$. To calibrate and simulate the SDE, we convert the SDE model terms into differentiable compute graphs (using the `diffrax` package) and calibrate using gradient descent (using the `optax` package)[27]. This reduces the reliance on more complex calibration

---

[3]https://github.com/rapidsai/cudf.

Figure 2: *Risk analyst agentic workflow.* The *risk analyst* uses a *builder* (dashed boxes) and *critic* (dotted boxes) framework for model discovery, in which the *builder* is tasked with implementation and simulation of the time series. Meanwhile, the *critic* calculates evaluation metrics on the simulated data and uses this to proposed new models to be implemented by the *builder*.

methods, scales to an arbitrary number of parameters, and allows for a simple loss functions calculated around either summary statistics or mean absolute error (MAE) between outputs. In this work, we calculate the loss, $L$, as the MAE between synthetic and historical price paths,

$$\mathcal{L}(\theta) = \mathrm{MAE}(S_{t;\theta}, S_{t;\theta}^*) = \frac{\Sigma_i^{N^{\mathrm{sims}}} |S_{t;i,\theta} - S_{t;i,\theta}^*|}{N^{\mathrm{sims}}}, \tag{2}$$

for $N^{\mathrm{sims}}$ simulation paths and where we look at the time average, comparing per timesteps across Monte Carlo realisations of the SDE. We note that we also considered function-based calibration which minimises, for example, the difference between moments of the time series but this resulted in worse results overall. However, these summary statistics are used when proposing models to test and calibrate during the model discovery loop, as outlined in Section 2.1.2. Note, that we have also investigated other statistical properties to evaluate the time series, i.e. the Fourier and signature transform, but found that these dramatically increased the computational time without offering material improvements.

### 2.1.2 Model Discovery

Agentic frameworks refer to the use of systems of LLMs and text-based sub-routines which are designed to complete a pre-determined task (usually complex and multi-step). These sub-routines include function calling, text summarisation, code debugging, and other tasks that LLMs have shown considerable capabilities at in the recent years [30–32]. Here, we introduce an agentic framework for the discovery of mathematical models of stochastic time series, in which programs for simulating the time series are proposed, implemented, assessed and scored by LLMs. We refer to this as a model-discovery loop in reference to the seminal work by George Box [26], and inspired by [24]. Our own model-discovery loop is set-out in Figure 2.

The process of model discovery follows a typical *builder-critic* feedback, where one set of tasks relates to the implementation of the data-generating model, with the aim that the model output is to be representative of historical observations. We refer to this set of tasks as the *builder*. Additionally, the synthetic data are then evaluated and compared to the historical data using a combination of accuracy metrics such as the moments of the time series. Comparisons between the distributions and these are then scored and used to generate model improvements. These tasks are collectively referred to as the *critic*.

Our aim is to have the *builder* and *critic* collaborate in an 'open-ended' environment [33]. However, there is a significant technical challenge in giving the LLM unbounded control over both the design

and implementation of a stochastic model. To balance structure and flexibility, we constrain the LLM's invention within an SDE framework, as set out in Section 2.1.1. While the framework provides mathematical structure, the LLM itself has full autonomy in suggesting and designing the SDE model. We give more details on both *builder* and *critic* below.

First, we consider the *builder*. Here, the primary goal is to convert a specific SDE equation into a working calibrated model that can generate appropriate historical data. The *builder* can be expressed as

$$\mathcal{B} : (m_n, c_n) \rightarrow F_n(t; \theta) = S_{t,n}, \tag{3}$$

where the candidate model, $m_n \in \mathcal{M}$, is taken from the space of proposed models, $\mathcal{M}$, and where $c_n \in \mathcal{C}$ is the context for choosing this model, for $n = 1, ..., N$ iterations of the model-discovery loop. We set the context as the history of previously proposed and tested models which we refer to as the *memory*. Additionally, $F_n(t; \theta)$ is the program implementing the model, $m_n$, given appropriate parameters, $\theta$ and which generates the stochastic price paths, $S_{t,n}$. The *builder* includes the following sub-tasks:

1. Implementing: Given the proposed model, $m_n$, the 'implementer' constructs a simulator, $F_n(t; \theta)$, which generates synthetic stochastic time series. We use a reference implementation as context for the implementer (using the `diffrax` package). We note that this reference implementation reproduces geometric Brownian motion, which is the first starting model used in the model-discovery loop, and the reference is static through subsequent iterations. The implementer includes both the initial construction of the program, and a parameter estimation step in which the model suggests an initial value, $\theta^0$, for the model calibration.

2. Simulating: Given the simulator, $F_n(t; \theta)$ and initial estimates for parameters, $\theta^0$, we then calibrate the model using the historical data, as outlined in Section 2.1.1, to determine parameters, $\theta$, and use this to generate Monte Carlo realisations of the stochastic time series, $S_{t,n}$, approximating the historical price paths.

Following the construction of the model, the synthetic data are then passed to a second sub-routine which we refer to as the *critic*, expressed by

$$\mathcal{R} : (S_{t,n}, c_n) \rightarrow c_{n+1}, m_{n+1}, \tag{4}$$

where $S_{t,n}$ is the set of Monte Carlo realisations of the price paths generated from the simulator; $c_n$ is the current agent context, which involves all previously tested models and their evaluation metrics, $c_{n+1}$ and $m_{n+1}$ are the updated context and the chosen next model to test. The *critic* includes the following tasks:

1. Testing: We calculate the evaluation metrics of the synthetic data, $e(S_t, S_t^*)$. These include summary statistics of the historical and simulated time series such as statistical moments (mean, variance, kurtosis, skew), tail metrics for the 95 percentile, the Hill estimator, Hurst exponent, jump intensity, autocorrelation function (ACF) metrics, growth rate, and the mean absolute deviation from the median [34, 35]. These statistics were chosen because they carry a large degree of information about the features of the generated time series and are known to be of relevance for financial time series[36, 37]. Alongside this, we consider direct tests between the distribution of returns of historical and synthetic time series for which we calculate the Kolmogorov-Smirnov test, Wasserstein distance, and the mean absolute percentage error [38–40]. Finally, we use an additional agent to construct a summary of the tested model, $m_n$, used to generate time series, $S_{t,n}$.

2. Scoring: The above evaluation metrics are then passed to a scoring agent, which provides an overall score based on both the loss of the metrics, $\mathcal{L}(\theta)$, and a novelty score, $K_n$. The latter, is
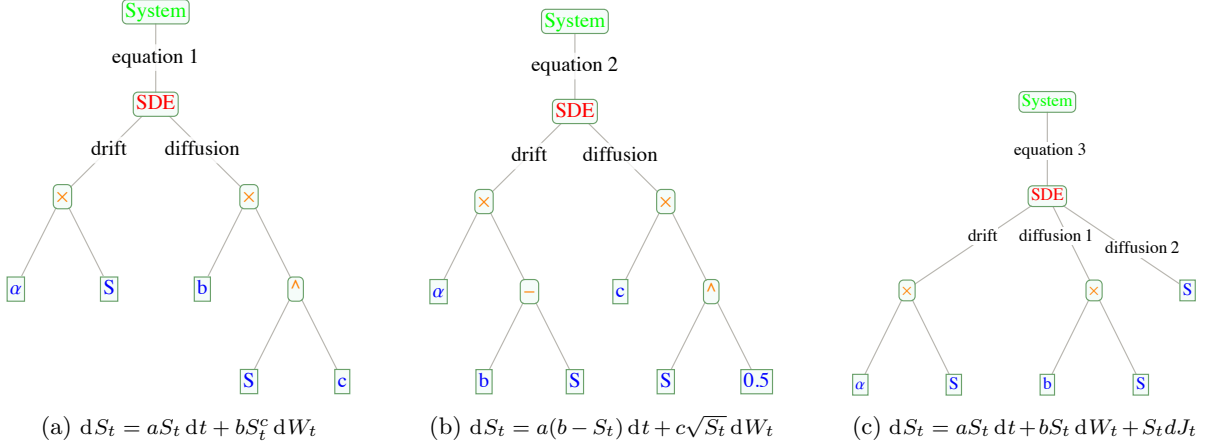
Figure 3: *Symbolic trees.* We convert the stochastic differential equation, used to simulate the financial time series, into a directed acyclic symbolic tree which we then use to calculate symbolic similarity and diversity scores. Here, we show (a) the CIR model, (b) the CEV model, and (c) the JD model, which are described in Section 3.1.

derived from an additional agent, using the model summary to determine a score between 1-100 on the novelty or 'interesting-ness' of the mathematical equation. This follows previous work that highlights that LLMs are potentially approaching human-like effectiveness at determining 'interesting-ness' due to the human biases implicit in their training data [32, 33, 41]. While we note that a simple scoring function can be estimated, which weights different components of the evaluation metrics to balance, for example, novelty versus fidelity, we consider a purely agentic method here to remove any reliance on tuning hyper-parameters and make a more principled comparison against LLMs. The model context is then updated to the model context, $c_{n+1} \in \mathcal{C}$, to determine the next model for testing.

3. Suggesting: Given the list of proposed models $\mathcal{M}$, their calibration loss, $\mathcal{L}(\theta)$, evaluation metrics, $e$, and novelty score, $K$, an LLM is then tasked with generating new candidate models that are expected to improve on previous iterations. This is based on the assumption that LLMs are able to approximate human creativity through their pre-training [42, 43]. These candidate models are then appended to the candidate model space.

4. Choosing: From the space of candidate models $\mathcal{M}$, the critic determines the next model to pass to the builder for implementation, $m_{n+1}$. This choice is based on the context $c_n$ of previously tested models.

The output of the *critic* loop is the next candidate model for the *builder* to implement and simulate. The model discovery process then repeats for $N$ iterations. Our model discovery loop involves repeated calls to LLMs which are tasked with specific agentic sub-routines. The orchestration of agents as either *builders* or *critics* is achieved using DSPy as the general framework for implementing agentic calls [28, 44]. Example outputs from the model-discovery loop agents are given in Appendix A.

### 2.1.3 Symbolic Similarity Score

In order to quantify the degree to which our model discovery loop is able to correctly identify a mathematical model, we construct a symbolic similarity score, $K_{\text{similarity}}$. Thus, we convert the proposed symbolic expression, describing model $m \in \mathcal{M}$, in to a directed acyclic graph in which each node represents the variable or parameter and edges represent the mathematical operation. We further decompose

the model into stochastic and deterministic components, representing the drift and diffusion terms of the underlying SDE. Figure 3 shows the symbolic trees for the three target equations discussed in Section 3.1.

Given a symbolic tree, we calculate graph-based statistics such as the mathematical complexity (the average number of paths within the tree) and the Weisfeiler-Lehman (WL) graph hash (a simple graph-based representation of the equation [45]). We then use the WL graph hash to calculate the similarity between two graphs, $K_{\text{similarity}}$, given by the following equation:

$$K_{\text{similarity}} = K_{\text{WL}}(G_i, G_j), \tag{5}$$

in which $G$ is the symbolic graph generated from the proposed and tested model, $m$, and $K_{\text{WL}}$ is the WL graph kernel calculated from WL graph hashes. The diversity of the proposed stochastic models, discovered by our *risk analyst* agents, $K_{\text{diversity}}$, is calculated as one minus the average pairwise similarity of all unique symbolic graph pairs, generated from the tested models, $m_n$, namely

$$K_{\text{diversity}} = 1 - \frac{1}{\binom{N}{2}} \sum_{1 \leq i < j \leq N} K_{\text{WL}}(G_i, G_j), \tag{6}$$

for $N$ models generated by the model discovery loop. Finally, we also calculate the overall complexity of the symbolic equation, a proxy for model complexity, $K_{\text{complexity}}$. This is defined as the average path length for the tree, in which larger equations with more terms and mathematical operations result in a higher complexity score. We next outline how the risk metrics are calculated from the final chosen model. These metrics are used by the *trader* agents to determine trading decisions, discussed in Section 2.2.

### 2.1.4 Model-informed Risk Metrics

The above model-discovery loop is used to determine a mathematical model for simulating the stochastic process that best matches a historical financial time series. Stochastic processes are widely used in financial trading applications, including option pricing, portfolio optimisation, and quantifying market risk [46, 47]. In this research study, we focus on market risk analysis, using the final model proposed by the model-discovery loop (Section 2.1.2) to calculate particular market risk statistics. Specifically, we calculate the value-at-risk (VaR), conditional value-at-risk (CVaR) and maximum drawdown (MDD), which are fundamental to risk management in financial markets [48], [49].

Given the generated model of the price path of an asset, denoted by $S_t$, we consider the loss to be $L_t = S_0 - S_t$, VaR is then a measure of these potential losses, defined as the quantile of the loss distribution for a given confidence level $(1 - \alpha)$,

$$\text{VaR}_\alpha = \inf\{l : P(L_t \leq l) \geq \alpha\}. \tag{7}$$

The CVaR, calculates the expected loss, given that we exceed VaR,

$$\text{CVaR}_\alpha = E[L_t | L_t > \text{VaR}_\alpha] \tag{8}$$

giving the magnitude of expected losses, rather than only the threshold.

Alongside VaR and CVaR, the MDD of a strategy over some period $[0, T]$ is commonly used to identify the relative maximum losses (peak-to-trough) incurred,

$$\text{MDD} = \max_{0 \leq t \leq T} \left[ \frac{\max_{0 \leq r \leq t} (S(r) - S(t))}{\max_{0 \leq r \leq t} S(r)} \right]. \tag{9}$$

These measures are fundamental in estimating risk exposure, but are known to underestimate inherent tail risk [50]. Given an appropriate model of market prices, Monte Carlo (MC) sampling can be

8

used to to generate $N^{\text{sims}}$ price paths from the SDE. These paths will be used to calculate both VaR and CVaR to approximate the tail risk.

Finally, we also use extreme value theory (EVT) to directly estimate the model- and market- specific tail risk. We consider the residuals, $\epsilon_t$, and the price paths to quantify model adequacy and to imply tail behaviour, respectively. For the residuals, given an appropriately calibrated SDE, we extract the standardized residuals as

$$\epsilon_t = \frac{S_{t+\Delta t} - S_t - f(S_t, t)}{g(S_t, t)\sqrt{\Delta t}} \tag{10}$$

for drift and diffusion terms, $f(S_t, t)$ and $g(S_t, t)$, respectively.

For a value of interest, $X$, we calculate the exceedances given a pre-defined threshold, $u$. Hence, we follow the peaks-over-threshold method and fit a generalized Pareto distribution (GPD) to the exceedances [51, 52],

$$P(X > x | X > u) = \left(1 + \xi \frac{x - u}{\gamma}\right)^{-1/\xi}, \tag{11}$$

in which $\xi$ is the shape parameter, for which we assume $\xi \neq 0$, and $\gamma$ is the scale parameter. The shape parameter gives estimate for the $k^{\text{th}}$ moment, given $\xi < 1/k$, and the scale parameter denotes the dispersion of these exceedences. By applying EVT to both residuals and the price paths, these parameters provide insights into the quality of the model fit to the historical market dynamics and the model-implied extreme behaviours.

When applying EVT to the simulation paths directly, we calculate the losses, $L_t$, and then fit the conditional excess above these losses using a threshold of 95% of the empirical loss distribution. For residuals, this calculates the extent to which our simulator has correctly reproduced the tail risk for the historical market price paths. Alternatively, for the losses of the simulated price paths, we obtain a better estimate for the tail risk than would be possible using historical data only, by sampling over Monte Carlo realisations.

We note that the above implementation assumes that the excess losses are independent which is typically not the case for financial time series. Further steps, such as through de-clustering, could be introduced to improve upon the risk analysis described above. However, we restrict our results to the above simplified implementation due to the challenges in automating more advanced risk analysis which we set as a topic for future research. We next set out how these measures are incorporated into trading strategies.

## 2.2 Trader Agents

Having determined a stochastic model of the time series, we then use the stochastic simulator to generate an array of price paths which are used to quantify market risk, and inform a trading strategy. In order to construct a trading strategy, we take an agentic approach and construct *trader* agents which make trading decisions based on model-informed risk and trend metrics, as well as market context taken from recent news headlines. The *trader* agents consist of a separate series of sub-routines, executed by an LLM. Here, we focus on equity markets, considering individual stocks. However, we note that our approach is broadly applicable to any asset class. We outline our methodology for combining market risk and price trends in a trading strategy below.

### 2.2.1 Trend Metrics

In order to give the LLM enhanced context for making trades, we use the relative strength index (RSI) to show the current trend of the price path [53]. RSI is a measure of signal momentum and is commonly
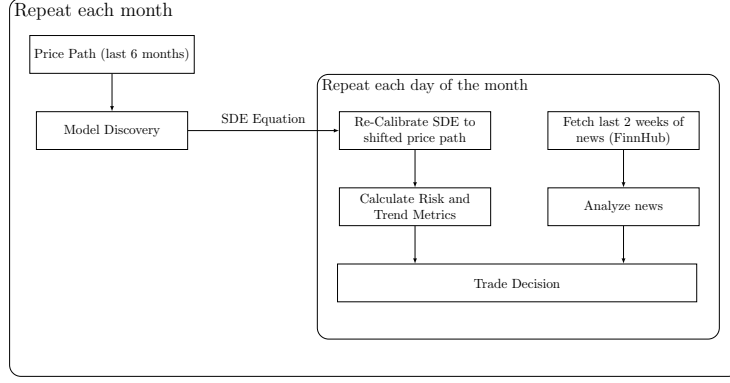
Figure 4: *Overview of the trading strategy pipeline.* Each month, the model discovery phase identifies a suitable SDE based on the past six months of price data. On each trading day, the SDE is recalibrated to the latest market conditions, recent news is fetched via the FinnHub API. Risk, trend metrics and sentiment are analyzed. These inputs are then used to inform a daily trade decision.

used for enhancing trading decisions by validating trends and trend reversals. It is defined as

$$\text{RSI} = 100 - \left( \frac{100}{1 + \frac{\text{Average Gain}}{\text{Average Loss}}} \right). \tag{12}$$

An RSI value below thirty usually indicates an oversold stock whereas a value above seventy indicates an overbought stock [53]. We set the lookback period for the loss and gain to fourteen days, as outlined in [53].

Given that the model discovered contains a drift term, it would be useful to exploit this to determine the trend of the signal. The drift term is separated from the SDE and the current price is passed to it. The standard SDE form is expressed in Equation 1. The drift term represents the expected rate of change of price, $S_t$. If the value of the drift is positive/negative then the price is expected to increase/decrease. This is passed to the LLM as a string for either 'positive' or 'negative'. We call this the Drift Polarity$_t$ term,

$$\text{Drift Polarity}_t = \begin{cases} \text{'Positive'}, & \text{if } f(S_t, t) \geq 0 \\ \text{'Negative'}, & \text{otherwise} \end{cases} \tag{13}$$

### 2.2.2 Trading Strategies

We have set out how LLMs can be used in an agentic framework for model discovery with a focus on the determination of a stochastic process. The generation of the SDE attempts to reproduce historical price dynamics. The overall class of SDEs is unconstrained and the individual model is determined by a *builder-critic* loop. We next set out how the outcome of this loop, the proposed SDE model, is used in a trading strategy.

To determine the role and impact of LLMs in trading applications, we consider a purely generative approach to trading decisions. To do so, we pass the calculated risk and trend metrics (VaR, CVaR, MDD, EVT metrics, RSI and Drift Polarity) to an LLM that is instructed to make trading decisions. We call this the *trader* agent and the framework for this is provided in Figure 4. LLMs are well-known for their ability to take insights from large context windows [54]. To take advantage of this we also provide context in the form of news articles to the trading agent and the methodology for this is outlined in Section 2.2.3.

At the beginning of each month, we run the model discovery loop and the best SDE (determined

10

by lowest loss) is chosen as the optimal model used for the next twenty days of trading. Here, we restrict the number of times we perform model discovery to monthly due to computational constraints. However, we found that the suggested models were relatively static over multiple consecutive trading days. Future work will look to quantify the consistency of the proposed models over time.

For each day of the month, we then recalibrate the SDE model to the price path of the last hundred trading days. The risk and trend metrics are then calculated from this calibrated SDE, as set out in Section 2.2.1. The summary of the news context is constructed by an LLM, which we refer to as the *news analyst*, described in Section 2.2.3. This summary is then passed to the *trader*, who makes the final decision to buy, sell, or hold. Example output from the *trader* is given in Appendix A.

The *trader* generates a trade signal, $T_t$, which can take one of three values: buy, sell, or hold. Based on this signal, the portfolio position, $P_t$, at time $t$ is updated according to the following rules:

$$P_t = \begin{cases} \frac{C_t}{S_t}, & \text{if } T_t = \text{buy and } C_t > 0 \\ 0, & \text{if } T_t = \text{sell} \\ P_{t-1}, & \text{if } T_t = \text{hold} \end{cases} \tag{14}$$

for $C_t$, the amount of cash at time $t$. $C_t$ is updated based on execution and transaction costs, such that

$$C_t = \begin{cases} C_{t-1} - P_t S_t - \text{TC}_t, & \text{if } T_t = \text{buy} \\ C_{t-1} + P_{t-1} S_t - \text{TC}_t, & \text{if } T_t = \text{sell} \\ C_{t-1}, & \text{if } T_t = \text{hold} \end{cases} \tag{15}$$

in which the transaction cost, $\text{TC}_t$, is determined by

$$\text{TC}_t = \kappa P_t S_t \tag{16}$$

for an assumed transaction charge, $\kappa$, given in basis points (bps). NYSE stocks are estimated to have an average charge of 8.8 bps of the total value of the trade, such that $\kappa = 8.8 \times 10^{-4}$[55].

Finally, we compare the LLM trading performance to a buy-and-hold strategy. While the buy-and-hold strategy is remarkably simple, it remains highly robust for capitalising on a bullish market, as we report in Section 3.2.

### 2.2.3 Trading Context

Over medium to long-term time horizons, a variety of other factors inform trading decisions ranging from macro-economic events, the performance of associated assets and companies, news sentiment, etc. To reproduce the effect of these factors on trading decisions, we pass market context to our *trader* agent.

Here, we use the `FinnHub` API[4] to gather context by retrieving news for the selected company within the specified date range. As outlined in Section 2.2, we consider trading over several months with trading decisions made daily. Given this, we use news context from the last five days. The news context provided to the *trader* is designed to give the agent access to the breadth of detail around current events. The `FinnHub` API summarises the news article into a short paragraph. This summary is key to keeping the context window short while still conveying the core message of the article. Each news item includes the timestamp of publication, the source of the article, the headline, a brief summary and the company that this item applies to.

We note that the movement of a company's stock price is typically highly dependent on the progress of other related companies. To account for this, we construct a *news analyst* agent that collects news around a specific company of interest by identifying other companies that may affect the price of the
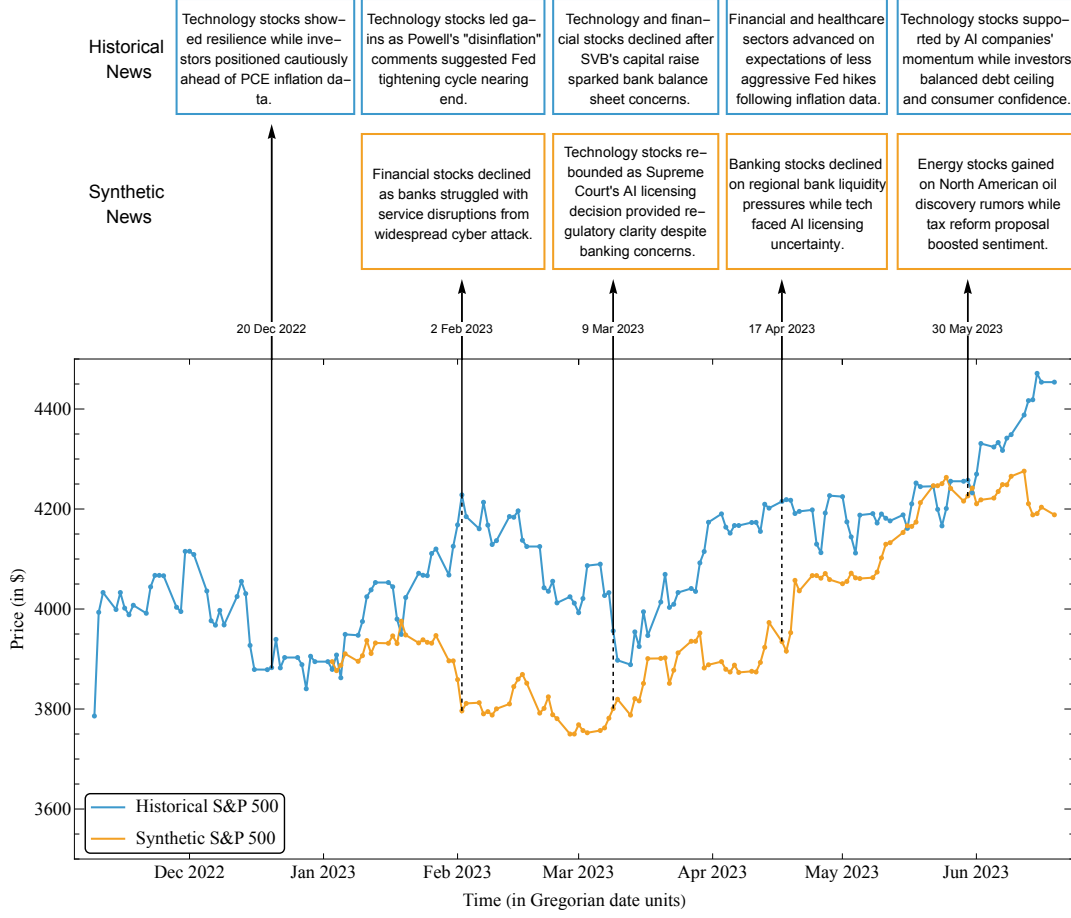
---

[4]Further information can be found here.

Figure 5: *Synthetic news generation.* To ensure that the tested LLMs do not have relevant context in their pretrained weights to improve trading results, we use the `Simudyne Horizon` simulator to generate synthetic price and news data. This provides a more advanced form of context-aware backtesting where trading strategies are tested on different plausible news cycles. In the Appendix B we give the full output of the generated news articles.

company in question. We then retrieve news for these related companies, with the assumption that this additional information of the wider sector/industry will improve trading decisions.

The *news analyst* agent determines the pros and cons to investing in this particular asset, according to the list of news items, and provides a summary of context-specific insights. Finally, these insights are passed to the *trader* to inform their trading decision. Example output from the *news analyst* is given in Appendix A.

### 2.2.4 Synthetic Trading Context

In the above description, we have set out a strategy of implementing and improving stochastic models for quantifying market risk using LLMs. We use these stochastic models to calculate risk and trend metrics that are passed to *trader* agents to determine trading (buy/sell/hold) decisions. However, when testing the performance of the *trader* agents, we could not preclude the possibility that historical events

had, under some circumstances, been accessible to the pre-trained LLMs, despite all testing periods being outside of the models knowledge cut-off dates.

To rule out this biasing we use `Simudyne Horizon`[5] to generate synthetic time series which have the statistical properties of historical price paths but relate to a set of synthetic market events. These events are described through news stories generated by LLMs, which blend real historical events with fictional ones by aligning recent price changes to plausible but fictional news stories such as cyber-attacks, alternative policy announcements, or other geopolitical/market events. This means that trading decisions informed by both price data and market context (such as when driven by sentiment analysis) can be effectively tested over a range of extreme scenarios without putting capital at risk. We show an example of the `Simudyne Horizon` context-aware backtesting functionality in Figure 5 and give examples of fictional headline summaries in Appendix B.

We use this as a testing environment to evaluate our *trader* agents, ensuring they cannot be influenced by data from their pre-training even when events are outside of their knowledge cut-off date. This bias may occur, for example, when events have a high likelihood of occurring due to past historical events, such as conservative estimates around rate decisions (i.e., the decision by central banks to typically make small changes to interest rates).

# 3    Results & Discussion

$\mathscr{D}$ ESPITE widespread attention to LLMs, there has been relatively little research on whether they can be used extensively for discovering financial models. Below, we outline our results in applying agents to both the discovery of stochastic models and the application of them to inform trading decisions. Firstly, we outline results comparing seven of the frontier LLMs that represent state-of-the-art performance at the time of writing for their ability to perform automated model discovery. Secondly, we evaluate the LLMs performance as risk-informed agentic traders, considering periods which we expect to be out of the LLMs training data. Lastly, we test the LLMs performance using a context-aware backtesting environment generated by `Simudyne Horizon`.

## 3.1    Discovering Stochastic Models

In order to quantify the degree to which our model discovery loop is effective at identifying reasonable candidate functions for risk analysis, we first test our framework on known SDEs. We consider three SDEs that are commonly used in finance and economics. These are the constant elasticity of variance (CEV) model, the Cox-Ingersoll-Ross (CIR) model, and a simple jump-diffusion (JD) model [56–58]. We give more details of each below. For each model, we use the same lowercase Roman alphabet for parameter variables, i.e. $a, b, c, \ldots$, for consistency with other models suggested by our *risk analyst* agent and to simplify the generation of the symbolic trees, outlined in Section 2.1.3. This means that the underlying interpretation of the parameters is model-dependent and requires human oversight in their application. We note that these parameters are then calibrated and the model discovery loop can suggest any number of parameters.

First, the CEV model is described by

$$\mathrm{d}S_t = aS_t\,\mathrm{d}t + bS_t^c\,\mathrm{d}W_t \tag{17}$$

in which $S_t$ is the stochastic variable, representing the price of the security over time, $a$ is the strength

---

[5] `Simudyne Horizon` is a hybrid market simulator combining traditional interpretable methods for simulating financial time series with agentic workflows. `Simudyne Horizon` generates synthetic time series and news data across capital markets, reproducing realistic price dynamics over days to years, while also providing interpretable estimates for the trading activity as grouped by trading styles.

| LLM | Loss ($\mathcal{L}(\theta)$) | $K_{\text{Score}}$ | $K_{\text{similarity}}(m_i, m^*)$ | Combined Score | $K_{\text{diversity}}$ | $K_{\text{complexity}}$ |
|---|---|---|---|---|---|---|
| Equation 1: $\mathrm{d}S_t = aS_t\,\mathrm{d}t + bS_t^c\,\mathrm{d}W_t$ | | | | | | |
| Deepseek: R1 | 1.0 (±1.5) | 52 (±8.4) | 0.44 (±0.03) | 1060 (±110) | <u>0.72</u> | 3.9 |
| Anthropic: Sonnet-3.7 | **0.040 (±0.043)** | **65 (±1.8)** | 0.40 (±0.069) | **1240 (±220)** | 0.67 | 4.3 |
| Anthropic: Sonnet-3.5v2 | 60 (±34) | 50 (±3.9) | <u>0.44 (±0.076)</u> | 980 (±300) | 0.7 | 4.0 |
| OpenAI: 4o-mini | 75 (±0.10) | 32 (±2.1) | 0.25 (±0.030) | 390 (±20) | 0.69 | <u>4.4</u> |
| OpenAI: o1-mini | 75 (±1.0) | 52 (±4.5) | 0.33 (±0.041) | 680 (±130) | **0.88** | 3.6 |
| OpenAI: o3-mini | 30 (±41) | <u>52 (±5.7)</u> | 0.40 (±0.10) | 880 (±180) | 0.68 | 3.4 |
| Meta: llama-3.3 | 60 (±34) | 36 (±4.6) | **0.45 (±0.017)** | 690 (±54) | 0.66 | **4.7** |
| Equation 2: $\mathrm{d}S_t = a(b - S_t)\,\mathrm{d}t + c\sqrt{S_t}\,\mathrm{d}W_t$ | | | | | | |
| Deepseek: R1 | 0.33 (±0.19) | <u>64 (±2.2)</u> | 0.35 (±0.060) | <u>980 (±100)</u> | 0.72 | 3.8 |
| Anthropic: Sonnet-3.7 | **0.00058 (±0.00048)** | **76 (±3.3)** | 0.37 (±0.042) | **1250 (±230)** | 0.67 | **4.4** |
| Anthropic: Sonnet-3.5v2 | <u>0.26 (±0.23)</u> | 56 (±3.1) | 0.41 (±0.054) | 970 (±30) | 0.67 | 3.9 |
| OpenAI: 4o-mini | <u>0.36 (±0.16)</u> | 51 (±2.2) | <u>0.37 (±0.042)</u> | 820 (±82) | <u>0.77</u> | 3.8 |
| OpenAI: o1-mini | 0.39 (±0.12) | 58 (±8.0) | 0.34 (±0.060) | 770 (±300) | **0.86** | 3.8 |
| OpenAI: o3-mini | 0.42 (±0.039) | 61 (±8.2) | **0.42 (±0.093)** | 970 (±170) | 0.59 | 3.2 |
| Meta: llama-3.3 | 0.40 (±0.097) | 54 (±2.2) | 0.27 (±0.097) | 600 (±110) | 0.48 | <u>4.1</u> |
| Equation 3: $\mathrm{d}S_t = aS_t\,\mathrm{d}t + bS_t\,\mathrm{d}W_t + S_t\,\mathrm{d}J_t$ | | | | | | |
| Deepseek: R1 | 0.21 (±0.01) | <u>64 (±1.3)</u> | 0.42 (±0.17) | 1180 (±280) | 0.7 | <u>3.9</u> |
| Anthropic: Sonnet-3.7 | **0.0064 (±0.010)** | **69 (±2.5)** | 0.42 (±0.18) | **1270 (±470)** | <u>0.72</u> | **4.2** |
| Anthropic: Sonnet-3.5v2 | <u>0.078 (±0.072)</u> | 63 (±1.3) | 0.47 (±0.14) | <u>1190 (±170)</u> | 0.66 | <u>3.9</u> |
| OpenAI: 4o-mini | 0.21 (±0.011) | 46 (±2.2) | **0.48 (±0.16)** | 720 (±87) | 0.7 | <u>3.9</u> |
| OpenAI: o1-mini | 0.21 (±0.00010) | 64 (±4.4) | 0.34 (±0.21) | 870 (±520) | **0.83** | <u>3.9</u> |
| OpenAI: o3-mini | 0.22 (±0.0048) | **69 (±4.2)** | 0.45 (±0.15) | **1270 (±210)** | 0.67 | 3.7 |
| Meta: llama-3.3 | 0.21 (±0.0016) | 49 (±5.3) | 0.33 (±0.23) | 800 (±500) | 0.38 | 3.4 |

Table 1: Performance metrics for three different equations, in which $m^*$ relates to each of the ground-truth equations (1-3) for each subsection of the table (**Bold**=best, <u>underlined</u>=second best). Note that for all the metrics apart from Loss, larger values equate to a higher performance.

of the drift term and $b$ is the elasticity parameter, denoting the strength between volatility and the asset price and $dW_t$ denotes Brownian noise. This SDE is widely used in modelling asset price dynamics, especially when the volatility is thought to depend non-linearly on the asset price [56].

The CIR model is described by

$$\mathrm{d}S_t = a(b - S_t)\,\mathrm{d}t + c\sqrt{S_t}\,\mathrm{d}W_t \tag{18}$$

in which $S_t$ is the stochastic variable, and $a, b, c$, are model parameters. We note that the CIR model is typically used to model interest rates, where $S_t$ would be replaced by, e.g. $r_t$ to denote the interest rate and $a, b, c$ denotes the mean-reversion speed, long-term mean, and volatility strength, respectively [57]. Here, we use $S_t$ for consistency with the other tested SDEs. We include this SDE as it is relatively simple and widely used, though we consider only equity price paths in later sections.

Finally, our jump-diffusion (JD) model is described by

$$\mathrm{d}S_t = aS_t\,\mathrm{d}t + bS_t\,\mathrm{d}W_t + S_t\,\mathrm{d}J_t \tag{19}$$

in which $a$ is the expected rate of return, $b$ is volatility, which we assume to be constant, and $\mathrm{d}J_t$ is a jump process, typically modelled as a Poisson process, $\mathrm{d}N_t$, with intensity parameter $\lambda$ [58]. We note that the jump term is often further expanded as $(Y - 1)\,\mathrm{d}N_t$, for random variable $Y$, representing the (log-normal) intensity of jumps. Here, we specify the ground-truth target equation as a simpler jump diffusion term but allow the model discovery process to identify models of increasing complexity.

We find that the majority of LLMs are able to identify candidate models that can recreate the target path generated by each of the above SDEs. This is shown in Table 1, in which most models have a reasonably low loss. The loss is calculated from the calibration process, given in Equation 7, representing path similarity. However, across all target SDEs, we find that `Sonnet 3.7` (Anthropic) consistently outperforms other LLMs, typically by several orders of magnitude. The code generated by all LLMs was typically able to run and simulate the observed time series. Additionally, we observed
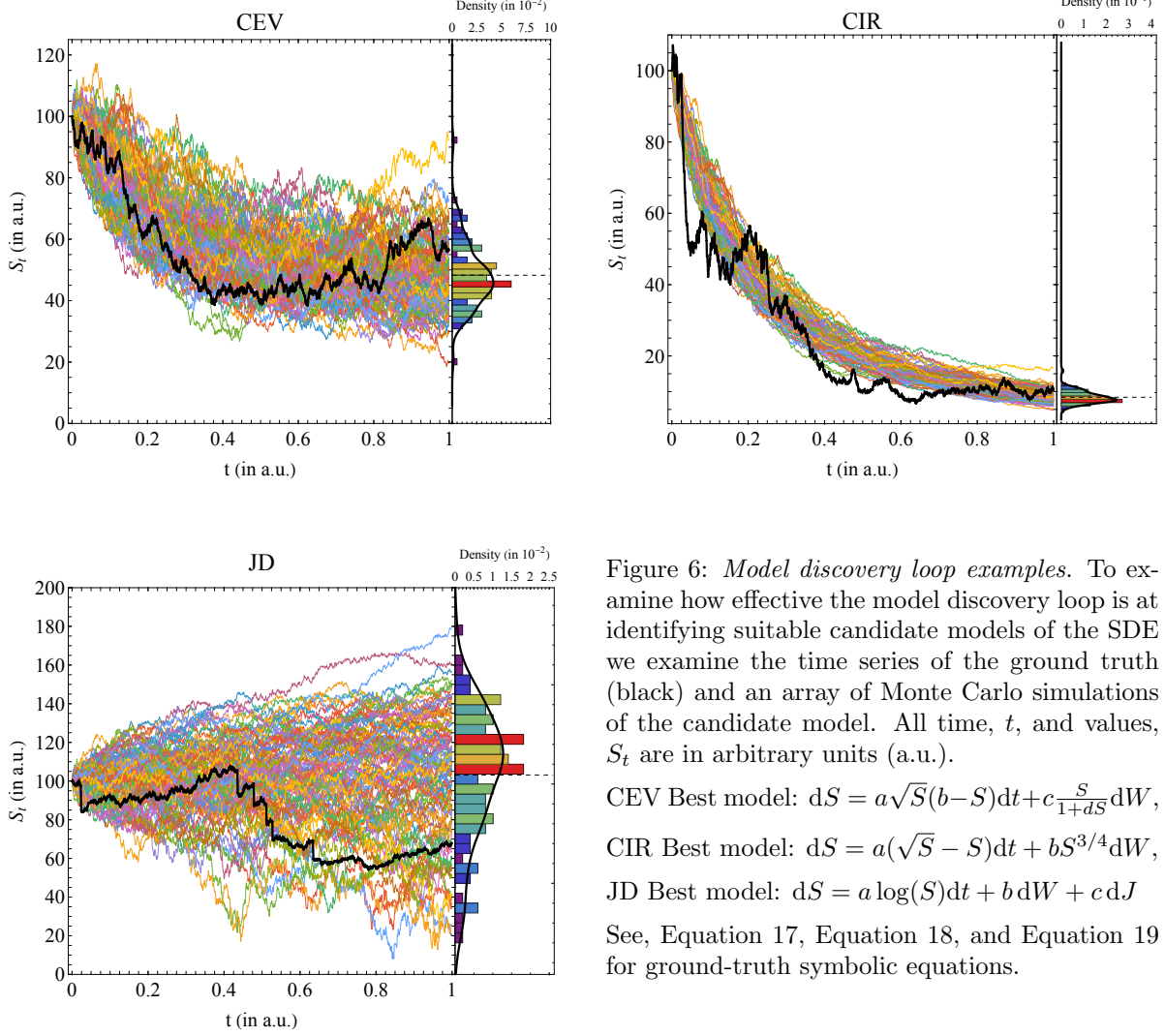
Figure 6: *Model discovery loop examples.* To examine how effective the model discovery loop is at identifying suitable candidate models of the SDE we examine the time series of the ground truth (black) and an array of Monte Carlo simulations of the candidate model. All time, $t$, and values, $S_t$ are in arbitrary units (a.u.).

CEV Best model: $\mathrm{d}S = a\sqrt{S}(b-S)\mathrm{d}t + c\frac{S}{1+dS}\mathrm{d}W$,

CIR Best model: $\mathrm{d}S = a(\sqrt{S} - S)\mathrm{d}t + bS^{3/4}\mathrm{d}W$,

JD Best model: $\mathrm{d}S = a\log(S)\mathrm{d}t + b\,\mathrm{d}W + c\,\mathrm{d}J$

See, Equation 17, Equation 18, and Equation 19 for ground-truth symbolic equations.

that different LLMs would sometimes choose and implement the same SDE with the same program, but these implementations would result in a different end result. In these cases, the differences in generating synthetic price paths is due to the parameter initialisation rather than the suggested SDE. This implies that `Sonnet 3.7` (the best model) is adept at not only program implementation but also initialising a stochastic model with a reasonable first estimate for the initial values.

While `Sonnet 3.7` was extremely good at both implementing and calibrating SDEs, when compared to other LLMs, the proposed models were often slightly different to the actual target SDE. This is shown in Figure 6. To quantify this, we use the symbolic similarity, described in Section 2.1.2, which uses the Weisfeiler-Lehman kernel to compare among symbolic trees, representing the underlying SDE. While `Sonnet 3.7` often suggests similar models to the ground-truth SDEs, $m^*$, we found that the most similar models were typically suggested by OpenAI series such as `o3-mini` and `4o-mini`. Additionally, `llama-3.3` (Meta) achieved very good model similarity for the first equation (CIR) but rarely suggested other models and often failed to correctly implement these models. In fact, when considering the diversity of the models suggested, $K_{\mathrm{diversity}}$, given by Equation 6, `o1-mini` (OpenAI) was considerably more creative and explorative in suggesting different model terms. These results highlight that LLMs will act in non-obvious and heterogenous ways within the same agentic framework, even when the

temperature is set equally (which we set to be 0.2 where possible in this research). This strengthens the need for robust testing when integrating LLMs into analysis workflows.

Finally, we note that we tested a mix of both reasoning and standard LLMs. In general, we find that reasoning LLMs, as is now well documented, typically perform better than their non-reasoning counterparts [3, 12]. This is true for `Sonnet 3.7` compared to `Sonnet 3.5v2` as well as `o3-mini` and `o1-mini` compared to `4o-mini`. In general, we find that LLMs that have reasoning capabilities have a 'reasoning bonus' when considering the combined score, which is a weighted combination of loss, LLM score and the similarity to the target equation. However, despite this 'reasoning bonus', we still find that non-reasoning LLMs can be highly performant such as the `Sonnet 3.5v2`, which was consistently within the top three best LLMs for these experiments.

## 3.2  Trading based on Stochastic Models

In this section we investigate the ability of LLMs to perform model discovery within realistic trading environments. We test the ability of each LLM to make trading decisions for the following equity symbols: *AAPL* (Apple), *NVDA* (NVIDIA), *F* (Ford) and *MSFT* (Microsoft). These companies were chosen due to the availability of news sources and for the interesting price dynamics over the chosen period. For example, NVIDIA experienced a significant fall in market price during this period, Microsoft and Ford had price paths with a similar start and end price, requiring a more active trading strategy in order to make a significant profit, and Apple had strong positive price movements.

For each asset, we repeat the model discovery loop fifteen times to find the most performant SDE model. This SDE is then used to generate the risk and trend metrics as described in Section 2.2.1, and these are then passed to an LLM (along with the news context) to make informed trading decisions. We evaluate this trading ability from 17$^{\text{th}}$ September 2024 to 24$^{\text{th}}$ April 2025. This was chosen to ensure that the trading period is outside of the LLMs training data [59]. We then repeat the model discovery process for the same LLMs as described in Section 3.1.

Each LLM is prompted to make a trading decision (buy, hold, or sell) given either solely news context from the last five days, $context^{\text{N}}$, or news context together with the risk and trend metrics, $context_{\text{M}}^{\text{N}}$, as generated by the model discovery loop. Trading decisions are then made on a daily basis over the trading period. Further details on the *trader* agent are given in Section 2.2. The trading strategy is then evaluated using three metrics:

1. Profit and Loss (PnL) (in $)

$$\text{PnL} = \text{Final Portfolio Value} - \text{Initial Portfolio Value}. \tag{20}$$

2. Sharpe ratio (SR)

$$SR = \frac{R_p - R_f}{\sigma_p} \tag{21}$$

   in which $R_p$ is the portfolio return, $R_f$ is the risk free return and $\sigma_p$ is the volatility of portfolio returns.

3. Maximum Drawdown (MDD) is given by Equation 9 and is a relative amount.

The initial portfolio value is $1000 and we use the three-month US treasury bill as our risk-free rate. For all symbols, we also compare our strategy against a simple buy-and-hold strategy.

Table 2 shows our results for the tested LLMs. We see an average SR of 0.88 when trading with just the news context. When including the model discovery loop to generate the risk and trend metrics, this average SR increases to 1.40. This shows that on average the model discovery loop enhances the SR of the trading strategies 37%.

| | R1 | Sonnet 3.7 | Sonnet 3.5v2 | 4o-mini | o1-mini | o3-mini | Llama 3.3 | B&H |
|---|---|---|---|---|---|---|---|---|
| | | | | *AAPL* (Apple) | | | | |
| PnL (N) | 215 | 61 | 100 | 347 | 166 | 293 | **431** | <u>372</u> |
| SR (N) | 2.39 | 0.40 | 0.84 | 3.70 | 1.50 | 2.67 | **4.56** | 3.53 |
| MDD (N) | 0.12 | 0.15 | 0.14 | **0.07** | 0.10 | 0.12 | **0.07** | 0.12 |
| PnL (N+M) | 230 | 337 | 75 | 328 | – | <u>384</u> | – | 372 |
| SR (N+M) | 2.49 | 3.27 | 0.60 | 3.05 | – | <u>3.87</u> | – | 3.53 |
| MDD (N+M) | 0.10 | <u>0.08</u> | 0.14 | 0.12 | – | **0.07** | – | 0.12 |
| | | | | *NVDA* (NVIDIA) | | | | |
| PnL (N) | 450 | 448 | 234 | 521 | 314 | 488 | 452 | 593 |
| SR (N) | 1.98 | 2.04 | 1.04 | 2.37 | 1.37 | 2.19 | 1.98 | 2.79 |
| MDD (N) | 0.27 | 0.27 | 0.25 | 0.27 | 0.31 | 0.28 | 0.28 | 0.27 |
| PnL (N+M) | 32 | <u>747</u> | – | 323 | -34 | 602 | **877** | 593 |
| SR (N+M) | 0.05 | <u>4.03</u> | – | 1.36 | -0.24 | 2.86 | **5.19** | 2.79 |
| MDD (N+M) | 0.39 | **0.23** | – | 0.28 | 0.38 | 0.27 | <u>0.24</u> | 0.27 |
| | | | | *MSFT* (Microsoft) | | | | |
| PnL (N) | 19 | 60 | 17 | 22 | 18 | 53 | 22 | 29 |
| SR (N) | -0.03 | 0.46 | -0.04 | 0.02 | -0.04 | 0.37 | 0.01 | 0.09 |
| MDD (N) | 0.15 | 0.11 | 0.16 | 0.12 | 0.12 | 0.15 | 0.15 | 0.15 |
| PnL (N+M) | <u>106</u> | **163** | <u>106</u> | -41 | 6 | – | – | 29 |
| SR (N+M) | <u>1.17</u> | **1.87** | 1.12 | -0.73 | -0.24 | – | – | 0.09 |
| MDD (N+M) | **0.07** | **0.07** | **0.07** | 0.15 | <u>0.08</u> | – | – | 0.15 |
| | | | | *F* (Ford) | | | | |
| PnL (N) | 6 | 5 | 78 | 63 | -46 | -140 | -130 | -146 |
| SR (N) | -0.13 | -2.26 | 0.92 | 0.67 | -0.40 | -0.83 | -0.94 | -0.84 |
| MDD (N) | 0.12 | **0.00** | 0.06 | 0.05 | 0.22 | 0.33 | 0.24 | 0.33 |
| PnL (N+M) | – | -306 | <u>179</u> | -145 | **180** | -44 | -72 | -146 |
| SR (N+M) | – | -1.86 | **2.96** | -0.98 | <u>1.82</u> | -0.37 | -0.57 | -0.84 |
| MDD (N+M) | – | 0.33 | <u>0.03</u> | 0.29 | 0.15 | 0.33 | 0.32 | 0.33 |

Table 2: *Trading results by symbol.* Profit and loss (PnL) is based on $1000 initial portfolio value, with units of dollars, while maximum drawdown (MDD) is given in relative terms. Sharpe ratio (SR) is unitless. N denotes news context (*context*$^N$) and N+M denotes news context with risk and trend metrics (*context*$_M^N$). B&H corresponds to a buy-and-hold strategy. Note that for PnL and SR, larger values equate to a higher performance, and for MDD a higher performance is depicted by a result closer to zero. **Bold**=best, <u>underlined</u>=second best and the dash (–) indicates that the model failed to complete the model discovery process.

This is especially true for `Llama 3.3` and `o3-mini`. For `Sonnet 3.7` and `Sonnet 3.5v2` the performance when comparing *context*$^N$ and *context*$_M^N$ is inconsistent. When the *context*$_M^N$ performance is better, however, this improvement is significant (roughly more than twice the PnL). As an example, when considering Microsoft, the PnL for `Sonnet 3.7` increases from 60 to 163 with the additional metrics context, similarly with `Sonnet 3.5v2` the PnL increases from 17 to 106, showing a large improvement.

For `4o-mini` and `o1-mini` the performance worsens with *context*$_M^N$ compared to *context*$^N$. These LLMs performed poorly at the model discovery step as shown in Table 1. This suggests that if the LLM is working with a poor model of the financial time series it will likely not be effective at producing profitable trading strategies. In addition to selecting suboptimal models, poor performance during model discovery is often caused by SDE models that fail to compile due to incorrect coding strategies. Rather than introducing complicated fallback solutions or attempting per-LLM prompt engineering, we tried to agentically debug each of these LLMs three times before marking this SDE as a failure. A loop failure is recorded when all fifteen SDE models selected during the model discovery process fail, we denote this as a dash (–), we show this in Table 2. We note that almost all LLMs showed at least one model failure during the model discovery process except for `Sonnet 3.7` and `4o-mini`, which we believe is due to their superior coding capabilities. The LLM with the highest failure rate is `Llama 3.3` as we discuss below.

Interestingly, `Llama 3.3` demonstrates inconsistent capabilities at model discovery. It suggests far more complex and novel symbolic models. However, this model complexity is not complimented by
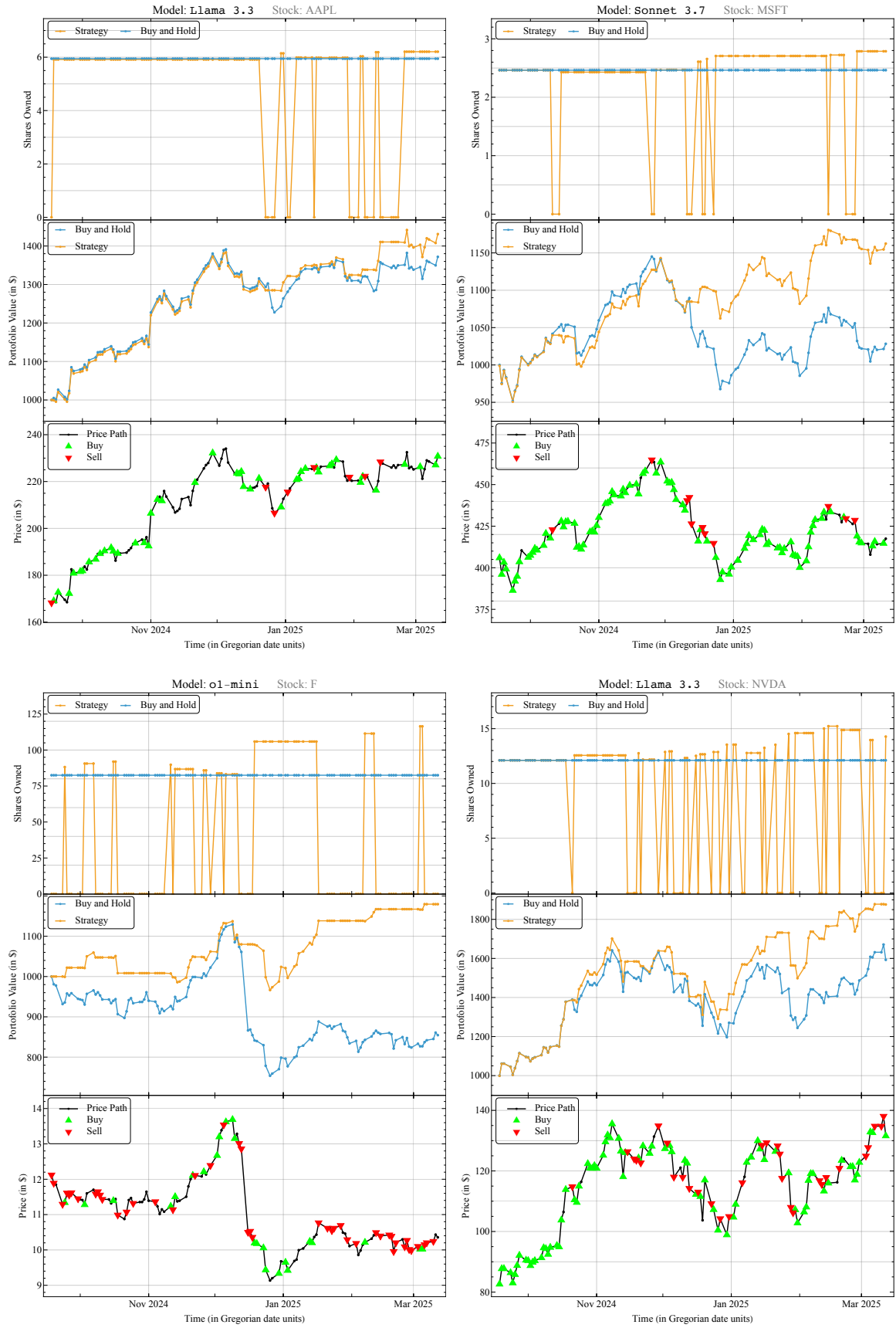
Figure 7: *Examples of trading results for each of the symbols considered.* Backtesting results from agentic trading strategies shown in orange for the period covering 17th September 2024 to 24th April 2025. We show results for `Llama 3.3` trading on *AAPL*, `Sonnet 3.7` trading on *MSFT*, `o1-mini` trading on *F*, `Llama 3.3` trading on *NVDA*. These represent the best results for each symbol, as shown in Table 2. Additionally, we show results from a buy-and-hold strategy in blue. Each panel shows trading decisions, the resulting portfolio output, and the number of shares owned during the trading period.

sufficient coding ability. This results in a worse performance at model discovery and we regularly could not get `Llama 3.3` to complete the full trading period.

In Table 1 we find that `Sonnet 3.7` performs the best at model discovery. However, we do not observe that this directly carries downstream to the results shown in Table 2. Poor model discovery typically results in poor trading. Conversely, it is not always the case that an LLM that builds a good candidate SDE (as measured by a low calibration loss, for example) will go on to make profitable trades. This is may be due to a number of factors, from poor initialisation of parameter values, changes in market regimes, unexpected news events, or simply due to the prompt-construction for the *trader* agent. Hence, the performance deficit between choosing an SDE model with a low calibration loss and making profitable trades will be due to our multi-component agentic framework (Figure 2). This framework requires LLMs to be equally good at suggesting SDE models to test, implementing these SDE models in code, analysing recent news articles and then balancing both model-informed risk statistics with news summaries to make a trading decision. This highlights the complexity inherent in building agentic trading systems. However, as discussed above, on average we find that the addition of a model discovery step improves trading results.

Next, we compare each LLM to the baseline strategy to get a better understanding of the performance of our agentic framework. We find that, when compared to a buy-and-hold strategy, the *trader* agents have mixed performance. For example, we find that for NVIDIA, no LLM is able to outperform the buy-and-hold baseline. This is likely due to the drastic increase in value of the NVIDIA stock over this period (roughly 47% returns), making outperforming a buy-and-hold very difficult. The buy-and-hold is outperformed by one LLM for Apple (`Llama 3.3`) and three LLMs for Microsoft (`Sonnet-3.7`, `Sonnet 3.5v2` and `R1`). Apple and Microsoft both saw returns of 28% and 3%, respectively. This shows the efficacy of the framework to generate notable returns beyond a bullish market. For Ford, all LLMs are able to outperform the baseline. This result is significant, Ford is the only symbol to see a decrease in value over the trading period. This means that in order to beat the baseline the *traders* needed to make decisions that took advantage of the short term price fluctuations in the market, something that is very difficult to do.

It is highly likely that different LLMs excel at different components within this framework, such as coding, sentiment analysis, and model discovery. Furthermore, we expect certain LLMs to perform better in non-obvious ways, such as the increased complexity of models suggested by `Llama 3.3`. Finally, the ability of LLMs appears to be linked to their reasoning ability. Models such as `R1`, `Sonnet 3.7`, and `o3-mini` are all reasoning LLMs and are often within the top two trading performances. However, this is not consistent across symbols. Future work is required to further evaluate these types of multi-agent LLM-based systems in complex tasks, as well as the application of smaller fine-tuned models.

## 3.3 Context-Aware Backtesting with `Simudyne Horizon`

In Section 3.2, we demonstrate that the combination of model discovery with news analysis can lead to profitable trading strategies. Additionally, we find that the distribution of performance across LLMs is non-obvious, i.e. reasoning LLMs may not correlate to trading performance. This is linked to the LLMs ability to perform sentiment analysis on recent events and model discovery for evaluating risk statistics. We assume that the overall performance of our *trader* agents will be influenced by the pre-training of the underlying LLMs, as we use a single LLM for all tasks within our agentic framework.

Assuming that all LLMs have been trained on a significant proportion of the internet up to their cut-off date, we cannot preclude the possibility of a biasing factor from their pre-training data affecting their trading performance in the equity markets. Additionally, it is not clear that we can entirely discount this even when using data explicitly outside of their training data knowledge cut-off as we have done in Section 3.1. This is due to the LLMs potentially using previous data in their pre-training to assign probabilities to unseen events which are implicit in the LLM inference output and additional to the

| | Deepseek R1 | Sonnet 3.7 | Sonnet 3.5v2 | 4o-mini | o1-mini | o3-mini | LLama 3.3 | B&H |
|---|---|---|---|---|---|---|---|---|
| PnL (N) | -62 | -12 | -22 | -71 | -13 | -113 | -28 | -99 |
| SR (N) | -0.77 | -0.68 | -1.23 | -1.06 | -0.44 | -1.46 | -0.85 | -1.09 |
| MDD (N) | 0.14 | <u>0.06</u> | **0.05** | 0.10 | 0.07 | 0.16 | 0.07 | 0.16 |
| PnL (N+M) | -88 | **53** | <u>45</u> | 1 | -66 | -55 | – | -99 |
| SR (N+M) | -1.70 | <u>0.39</u> | **0.47** | -0.21 | -1.07 | -1.13 | – | -1.09 |
| MDD (N+M) | 0.11 | <u>0.06</u> | <u>0.06</u> | 0.08 | 0.09 | 0.12 | – | 0.16 |

Table 3: *Trading results using the* `Simudyne Horizon` *backtesting environment.* Profit and loss (PnL) is based on \$1000 initial portfolio value, with units of dollars, while maximum drawdown (MDD) is given in relative terms. Sharpe ratio (SR) is unitless. N denotes news context ($context^{\mathrm{N}}$) and N+M denotes news context with risk and trend metrics ($context^{\mathrm{N}}_{\mathrm{M}}$). B&H corresponds to the buy-and-hold strategy. Note that for PnL and SR, larger values equate to a higher performance, and for MDD a higher performance is depicted by the value closest to zero. **Bold**=best, <u>underlined</u>=second best and the dash (–) indicates that the model failed to complete the model discovery process.

news summaries provided in $context^{\mathrm{N}}_{\mathrm{M}}$. Moreover, for those LLMs that do not expose their reasoning (via their API), we cannot be sure that system prompts that are unseen have been designed to augment their responses in non-obvious ways.

To address this, we use the `Simudyne Horizon` platform to generate an alternative price path and synthetic news cycle for the S&P 500 covering the 6 month trading period from 1$^{\mathrm{st}}$ January 2023 to 1$^{\mathrm{st}}$ July 2023. Illustratively, this includes fictional events such as the blocking of a major pharmaceutical merger, a major cyber attack on financial infrastructure, and the introduction of AI licensing by the US supreme court. Results from running the same trading agent and buy-and-hold scenario in the `Simudyne Horizon` backtesting environment are given in Table 3.

We find that most LLMs are able to outperform the simple buy-and-hold strategy by using news stories only, $context^{\mathrm{N}}$. The trading results are further improved when combined with a model discovery loop, as in Section 3.2. We see an average 22% increase in SR when comparing LLMs using $context^{\mathrm{N}}$ to $context^{\mathrm{N}}_{\mathrm{M}}$. The models `4o-mini`, `o3-mini`, `Sonnet 3.5v2`, and `Sonnet 3.7` all improved with the additional risk and trend metrics, however, both `R1` and `o1-mini` resulted in worse performance. As in Section 3.1, we find that `Llama 3.3` is unable to complete the model discovery due to the complexity of models that were suggested. Additionally, we find that `Sonnet 3.5v2`, and `Sonnet 3.7` and `4o-mini` resulted in profitable trading decisions across the synthetic backtest. Moreover, these results highlight that model-informed analysis improves agentic trading decisions and that these results are robust to potential pre-training biases. However, as in Section 3.2, the performance of the agentic *trader* agent is connected to the underlying LLM in a non-obvious way, such that reasoning LLMs do not necessarily result in improved trading performance. Our results highlight both the potential and fragility of agentic trading workflows.

# 4   Conclusion

$\mathscr{E}$VIDENCE in this research highlights that LLMs can be deployed in a novel agentic trading framework, leading to improved trading performance. Our research posed the title question on whether to trade or not to trade. We believe that this is no longer a question posed blindly to the market, nor answered solely by fundamental analysis, momentum, or gut instinct. In our research, we have shown that LLMs, embedded within an agentic framework, can participate meaningfully in the construction of explanatory models of financial time series. These models simulate, rather than simply predict, the uncertainty inherent in market dynamics. Yet perhaps the most significant contribution lies not in their performance gains, but in the semi-automated construction of an epistemic loop where discovery, simulation, and judgement co-evolve.

By integrating symbolic model-informed context into trading agents, capable of digesting both price

signals and narrative context, we have demonstrated a methodology for LLMs to make trading decisions. This method allows for model-aware risk quantification, in which the model development is partially automated by LLMs. This supports other recent works demonstrating that LLMs can be deployed to automate research and development tasks [60, 61]. We believe that this is the first time that this approach, LLM-based model development, has been applied to quantitative risk analysis for financial markets. Specifically, we demonstrate that our method can accurately determine the form of a target SDE, as outlined in Section 3.1.

Our model discovery loop was used to estimate the market risk to enhance trading strategies. Here, we observed that model-informed risk and trend metrics improved trading decisions, and that this represented a considerable improvement against a simple buy-and-hold strategy. However, we believe future work should consider the implications of such systems in the context of black swan events. In this research, we have found a distinct difference in the capabilities of the tested LLMs and expect that the specific skill-set of individual LLMs may be more effectively harnessed towards a common goal by applying them to LLM-specific sub-tasks. Future work will investigate the advantage of using different LLMs for individual tasks within a comprehensive agentic trading framework.

This study represents a first step towards the introduction of automated LLM-based reasoning agents in financial analysis and trading systems. We allow for informed and explainable decisions through the construction of interpretable programs representing simulators of financial time series. This provides clarity on decision-making. Most importantly, we highlight a considerable opportunity for profitable strategies. Yet, at the same time, we observe plausible risks inherent to the current generation of LLMs if used in a financial multi-agent framework. This is demonstrated by the LLMs inconsistent and varied performance in various steps of our agentic workflow. We expect that both opportunities and risks will compound in the future. As agentic systems mature, the question will not be to trade, but how to trade wisely — being mindful that the future is not predicted, but emerges through the recursive intelligence of models that learn to reason, simulate, and decide.

# 5 Acknowledgements

# References

[1] Liyang Wang, Yu Cheng, Ao Xiang, Jingyu Zhang, and Haowei Yang. "Application of natural language processing in financial risk detection". In: *arXiv preprint arXiv:2406.09765* (2024).

[2] Yuqi Nie et al. "A survey of large language models for financial applications: Progress, prospects and challenges". In: *arXiv preprint arXiv:2406.11903* (2024).

[3] Mohamed Amine Ferrag, Norbert Tihanyi, and Merouane Debbah. "From llm reasoning to autonomous ai agents: A comprehensive review". In: *arXiv preprint arXiv:2504.19678* (2025).

[4] Yangyang Yu et al. "FinMem: A performance-enhanced LLM trading agent with layered memory and character design". In: *Proceedings of the AAAI Symposium Series*. Vol. 3. 1. 2024, pp. 595–597.

[5] Han Ding, Yinheng Li, Junhao Wang, and Hang Chen. "Large language model agent in financial trading: A survey". In: *arXiv preprint arXiv:2408.06361* (2024).

[6] Alejandro Lopez-Lira. "Can Large Language Models Trade? Testing Financial Theories with LLM Agents in Market Simulations". In: *arXiv preprint arXiv:2504.10789* (2025).

[7] Huaqin Zhao et al. "Revolutionizing finance with llms: An overview of applications and insights". In: *arXiv preprint arXiv:2401.11641* (2024).

[8] Yuqi Nie et al. "A survey of large language models for financial applications: Progress, prospects and challenges". In: *arXiv preprint arXiv:2406.11903* (2024).

[9] Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. "Large language models in finance: A survey". In: *Proceedings of the fourth ACM international conference on AI in finance*. 2023, pp. 374–382.

[10] Bokai Cao et al. "From Deep Learning to LLMs: A survey of AI in Quantitative Investment". In: *arXiv preprint arXiv:2503.21422* (2025).

[11] Aske Plaat et al. "Agentic Large Language Models, a survey". In: *arXiv preprint arXiv:2503.23037* (2025).

[12] Zixuan Ke et al. "A Survey of Frontiers in LLM Reasoning: Inference Scaling, Learning to Reason, and Agentic Systems". In: *arXiv preprint arXiv:2504.09037* (2025).

[13] Aaron Jaech et al. "Openai o1 system card". In: *arXiv preprint arXiv:2412.16720* (2024).

[14] Daya Guo et al. "Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning". In: *arXiv preprint arXiv:2501.12948* (2025).

[15] Parshin Shojaee et al. "The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity". In: *arXiv preprint arXiv:2506.06941* (2025).

[17] Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. "tau-bench: A Benchmark for Tool-Agent-User Interaction in Real-World Domains". In: *arXiv preprint arXiv:2406.12045* (2024).

[18] Yijia Xiao, Edward Sun, Di Luo, and Wei Wang. "TradingAgents: Multi-Agents LLM Financial Trading Framework". In: *arXiv preprint arXiv:2412.20138* (2024).

[19] Saizhuo Wang et al. "Alpha-gpt: Human-ai interactive alpha mining for quantitative investment". In: *arXiv preprint arXiv:2308.00016* (2023).

[20] Zhizhuo Kou et al. "Automate strategy finding with llm in quant investment". In: *arXiv preprint arXiv:2409.06289* (2024).

[21] Qianggang Ding, Haochen Shi, Jiadong Guo, and Bang Liu. "Tradexpert: Revolutionizing trading with mixture of expert llms". In: *arXiv preprint arXiv:2411.00782* (2024).

[22] Yang Li, Yangyang Yu, Haohang Li, Zhi Chen, and Khaldoun Khashanah. "Tradinggpt: Multi-agent system with layered memory and distinct characters for enhanced financial trading performance". In: *arXiv preprint arXiv:2309.03736* (2023).

[23] Hoyoung Lee, Youngsoo Choi, and Yuhee Kwon. "Quantifying Qualitative Insights: Leveraging LLMs to Market Predict". In: *arXiv preprint arXiv:2411.08404* (2024).

[24] Michael Y Li, Emily B Fox, and Noah D Goodman. "Automated statistical model discovery with language models". In: *arXiv preprint arXiv:2402.17879* (2024).

[25] Namid R Stillman and Rory Baggott. "Neuro-Symbolic Traders: Assessing the Wisdom of AI Crowds in Markets". In: *arXiv preprint arXiv:2410.14587* (2024).

[26] G. E. P. Box and William G. Hunter. "A Useful Method for Model-Building". In: *Technometrics* 4.3 (1962), pp. 301–318. ISSN: 00401706.

[27] Patrick Kidger. "On neural differential equations". In: *arXiv preprint arXiv:2202.02435* (2022).

[28] Omar Khattab et al. "DSPy: Compiling Declarative Language Model Calls into Self-Improving Pipelines". In: 2024.

[29] Anna Orlik and Laura Veldkamp. "Understanding uncertainty shocks and the role of black swans". In: *Journal of Economic Theory* 222 (2024), p. 105905. ISSN: 0022-0531.

[30] Lochan Basyal and Mihir Sanghvi. "Text summarization using large language models: a comparative study of mpt-7b-instruct, falcon-7b-instruct, and openai chat-gpt models". In: *arXiv preprint arXiv:2310.10449* (2023).

[31] Nan Jiang et al. "Ledex: Training LLMs to better self-debug and explain code". In: *Advances in Neural Information Processing Systems* 37 (2024), pp. 35517–35543.

[32] Maxence Faldor, Jenny Zhang, Antoine Cully, and Jeff Clune. "OMNI-EPIC: Open-endedness via Models of human Notions of Interestingness with Environments Programmed in Code". In: *arXiv preprint arXiv:2405.15568* (2024).

[33] Edward Hughes et al. "Open-endedness is essential for artificial superhuman intelligence". In: *arXiv preprint arXiv:2406.04268* (2024).

[34] Ruey S Tsay. *Analysis of financial time series*. John wiley & sons, 2005.

[35] Rama Cont and Peter Tankov. *Financial modelling with jump processes*. Chapman and Hall/CRC, 2003.

[36] Sara A Safari, Maximilian Janisch, and Thomas Lehéricy. "International Financial Markets Through 150 Years: Evaluating Stylized Facts". In: *arXiv preprint arXiv:2504.08611* (2025).

[37] Ethan Ratliff-Crain, Colin M Van Oort, James Bagrow, Matthew TK Koehler, and Brian F Tivnan. "Revisiting stylized facts for modern stock markets". In: *2023 IEEE International Conference on Big Data (BigData)*. IEEE. 2023, pp. 1814–1823.

[38] Rob J Hyndman and George Athanasopoulos. *Forecasting: principles and practice*. OTexts, 2018.

[39] Cédric Villani et al. *Optimal transport: old and new*. Vol. 338. Springer, 2008.

[40] Myles Hollander, Douglas A Wolfe, and Eric Chicken. *Nonparametric statistical methods*. John Wiley & Sons, 2013.

[41] Cheng-Han Chiang and Hung-yi Lee. "Can Large Language Models Be an Alternative to Human Evaluations?" In: *Annual Meeting of the Association for Computational Linguistics*. 2023.

[42] Harsh Kumar, Jonathan Vincentius, Ewan Jordan, and Ashton Anderson. "Human Creativity in the Age of LLMs: Randomized Experiments on Divergent and Convergent Thinking". In: *arXiv preprint arXiv:2410.03703* (2024).

[43] Yufei Tian et al. "MacGyver: Are Large Language Models Creative Problem Solvers?" In: *arXiv preprint arXiv:2311.09682* (2023).

[44] Omar Khattab et al. "Demonstrate-Search-Predict: Composing Retrieval and Language Models for Knowledge-Intensive NLP". In: *arXiv preprint arXiv:2212.14024* (2022).

[45] Nino Shervashidze, Pascal Schweitzer, Erik Jan Van Leeuwen, Kurt Mehlhorn, and Karsten M Borgwardt. "Weisfeiler-lehman graph kernels." In: *Journal of Machine Learning Research* 12.9 (2011).

[46] Tomasz Rolski, Hanspeter Schmidli, Volker Schmidt, and Jozef L Teugels. *Stochastic processes for insurance and finance*. John Wiley & Sons, 2009.

[47] Masaaki Kijima. *Stochastic processes with applications to finance*. Chapman and Hall/CRC, 2002.

[48] Hongzhong Zhang, Tim Leung, and Olympia Hadjiliadis. "Stochastic modeling and fair valuation of drawdown insurance". In: *Insurance: Mathematics and Economics* 53.3 (Nov. 2013), pp. 840–850. ISSN: 0167-6687.

[49] R. Tyrrell Rockafellar and Stanislav Uryasev. "Optimization of conditional value-at risk". In: *Journal of Risk* 3 (2000), pp. 21–41.

[50]  Paul Embrechts, Rdiger Frey, and Alexander McNeil. *Quantitative Risk Management*. 2011.

[51]  Stuart Coles, Joanna Bawa, Lesley Trenner, and Pat Dorazio. *An introduction to statistical modeling of extreme values*. Vol. 208. Springer, 2001.

[52]  Alexander J McNeil, Rüdiger Frey, and Paul Embrechts. *Quantitative risk management: concepts, techniques and tools-revised edition*. Princeton university press, 2015.

[53]  Pat Tong Chio. "A comparative study of the MACD-base trading strategies: evidence from the US stock market". In: *arXiv preprint arXiv:2206.12282* (2022).

[54]  Yiran Ding et al. "Longrope: Extending llm context window beyond 2 million tokens". In: *arXiv preprint arXiv:2402.13753* (2024).

[55]  Robert F. Engle, R. Ferstenberg, and Jeffrey R. Russell. "Measuring and Modeling Execution Cost and Risk". In: *IO: Firm Structure* (2006).

[56]  Stan Beckers. "The constant elasticity of variance model and its implications for option pricing". In: *the Journal of Finance* 35.3 (1980), pp. 661–673.

[57]  Ludger Overbeck and Tobias Ryden. "Estimation in the cox-ingersoll-ross model". In: *Econometric Theory* 13.3 (1997), pp. 430–461.

[58]  Gerald HL Cheang and Carl Chiarella. "A modern view on Merton's jump-diffusion model". In: *Stochastic Processes, Finance and Control: A Festschrift in Honor of Robert J Elliott*. World Scientific, 2012, pp. 217–234.

[59]  Hao Wang. *LLM Knowledge Cutoff Dates*. GitHub repository. 2024. URL: https://github.com/HaoooWang/llm-knowledge-cutoff-dates.

[60]  Chris Lu et al. "The ai scientist: Towards fully automated open-ended scientific discovery". In: *arXiv preprint arXiv:2408.06292* (2024).

[61]  Juraj Gottweis et al. "Towards an AI co-scientist". In: *arXiv preprint arXiv:2502.18864* (2025).

# LLM Co-authored References

[16]  C Opus and A Lawsen. "The Illusion of the Illusion of Thinking". In: *arXiv preprint arXiv:2:2506.09250* (2025).

# Appendix

## A    Example Agent Output

### Summarising Agent

This is a non-linear SDE with interesting properties: Looking at the model code, we can see:

- The drift term is a*sin(S/b) where S is the state variable

- The diffusion term is c*S where c is a constant

- This forms a stochastic differential equation (SDE)

- The equation has a nonlinear sinusoidal drift term which is interesting

- The diffusion term is proportional to S which is standard

### Scorer Agent

Let's analyze this model in two aspects:

1. Novelty (50% of score):

    - The model uses a sine function which is a novel approach for SDEs
    - The structure dS = a*sin(S/b)dt + cSdW is relatively simple but introduces non-linear behavior
    - The sine function could capture cyclical patterns in the data
    - Score for novelty: 35/50 (good originality but maintains reasonable simplicity)

2. Metrics Accuracy (50% of score):

    - Delta metrics show good alignment in mean (0.001) and jump intensity (0.005)
    - Volatility is reasonably close (delta of 0.278)
    - ACF values are fairly close to target (deltas around 0.02-0.04)
    - Kurtosis and skewness show moderate differences
    - The range and distribution metrics (min, max, mean) show some deviation
    - Kolmogorov-Smirnov test value of 0.276 indicates moderate fit
    - Score for metrics: 32/50

Total score calculation: 35 (novelty) + 32 (metrics) = 67/100

## Model Suggestor Agent

- $dS = a(b - S)dt + c\sqrt{S}dW$

- $dS = aS(b - \ln(S))dt + cSdW$

- $dS = aS^b dt + cdW$

- $dS = kS(1 - S)dt + mS^{1/2}dW$

- $dS = k(S/(1 + S^2))dt + mSdW$

- $dS = a\tanh(bS)dt + cSdW$

- $dS = k(1 - \exp(-mS))dt + p\sqrt{S}dW$

- $dS = a/(1 + \exp(-bS))dt + c(1 + S)dW$

- $dS = k\sin(mS)dt + p\cos(S)dW$

- $dS = aS(1 - \ln(S/b))dt + cS^{1/3}dW$

- $dS = kS/(1 + S^3)dt + m\exp(-pS)dW$

- $dS = a\sin(S/b)dt + cSdW$

- $dS = kS/(1 + S^2)dt + mS/(1 + S)dW$

- $dS = k(1 - S^{1/3})dt + mS^{1/3}dW$

- $dS = aS\log(1 + bS)dt + c/(1 + S)dW$

- $dS = k\arctan(mS)dt + pSdW$

- $dS = a\exp(-bS)dt + c \cdot S/(1 + S^2)dW$

- $dS = a\sin(bS)/(1 + S^2)dt + cSdW$

- $dS = k \cdot S\log(1 + mS)dt + p \cdot S\log(1/(1 + S))dW$

- $dS = k\log(1 + S)/(1 + S^2)dt + m \cdot S/(1 + \log(S))dW$

- $dS = a \cdot S^{1/2}(1 - \exp(-bS))dt + c(1 + S^{1/3})dW$

## News Gatherer Agent

Based on the news articles, there's significant focus on DeepSeek, a Chinese AI startup that's causing market disruption in the AI sector. The news suggests this is affecting various tech companies, particularly those involved in AI and semiconductors. To make a comprehensive analysis of NVDA's current state, we should look at:

1. Direct competitors in the GPU/AI chip space (AMD, INTC)

2. Key semiconductor manufacturing partners (TSM)

3. Major AI technology companies affected by DeepSeek developments (MSFT, GOOGL)

4. Companies involved in AI infrastructure (ORCL)

These companies will provide a complete picture of the competitive landscape, supply chain, and broader AI market dynamics affecting NVDA.

## News Analyst Agent

The emergence of DeepSeek and other competitive threats suggests NVIDIA's dominant position in AI chips could face increasing pressure. The company's high margins and premium pricing strategy may become harder to maintain as more efficient, lower-cost alternatives emerge. The concentration of value in AI-related growth also creates vulnerability if AI spending patterns shift or if geopolitical tensions impact the global tech supply chain.

While NVIDIA remains a leader in its space, these factors combine to create meaningful downside risks that investors should carefully consider. Despite recent market volatility and concerns about competition from Chinese AI firms like DeepSeek, NVDA maintains strong fundamentals and growth prospects. The company's involvement in Project Stargate, continued innovation in AI hardware, and expansion into new markets like healthcare suggest sustained growth potential. Their dominant market position, strong financial performance, and analyst confidence indicate NVDA remains well-positioned to capitalize on the expanding AI market.

Recent dips in stock price may present buying opportunities for long-term investors, particularly given the company's track record of successful innovation and market leadership in AI technology.

Date: 2025-01-30

After analyzing the news and analyst thoughts, I recommend a SELL decision for NVIDIA stock for the following reasons:

1. **Valuation Concerns**: Multiple sources indicate AI chip stocks, including NVIDIA, are trading at "patently absurd" valuations with triple-digit P/E ratios that appear disconnected from fundamentals, suggesting an overvalued position.

2. **Emerging Competitive Threats**: The development of China's DeepSeek AI tool represents a tangible competitive threat that has already caused an 8% pre-market drop in NVIDIA's stock, indicating the market's sensitivity to competitive pressures in the AI chip space.

3. **AI Spending Bubble Risk**: There are substantial concerns about an AI spending bubble, with analysts warning of an "inevitable market correction" for AI tech stocks. This suggests NVIDIA's current price may be inflated by unsustainable investment patterns.

4. **Market Speculation Over Fundamentals**: Multiple articles indicate that NVIDIA and other AI stocks are trading based on speculation (particularly around "Project Stargate") rather than concrete financial performance, creating vulnerability to sentiment shifts.

5. **Short-Term Negative Catalysts**: The "DeepSeek problem" specifically identified for NVIDIA represents an immediate concern that could continue to pressure the stock in the near term.

While NVIDIA maintains a strong position in the AI market with potential long-term government support through initiatives like Project Stargate, the current valuation appears to have priced in overly optimistic scenarios while underweighting significant risks. The combination of excessive valuation, emerging competitive threats, and bubble concerns creates an unfavorable risk-reward profile for the immediate term, making a short position (SELL) the appropriate trading decision for daily position evaluation.

Decision: **SELL**

Date: 2025-01-30

Model: $dS = a(b - S)dt + c\sqrt{S}\,dW$

After analyzing the risk metrics and analyst thoughts, I'm recommending a HOLD position on NVIDIA for the following reasons:

1. **Mixed Risk Profile**:

   - The risk metrics show some concerning elements: The $EVT_{VaR}$ (20.79) and $EVT_{CVaR}$ (25.50) are significantly higher than the standard VaR (9.90) and CVaR (16.08), indicating that extreme negative movements could be worse than standard risk measures suggest.

   - The Maximum Drawdown (MDD) of 52.70% is substantial and indicates historical vulnerability to major corrections.

   - However, the negative EVT shape parameters for both losses (-0.15) and residuals (-0.13) suggest bounded tail distributions, which is somewhat reassuring.

2. **Neutral Technical Indicators**:

   - The RSI at 52.21 is in neutral territory, neither overbought nor oversold.

   - The positive Drift Score suggests some upward momentum in the price action.

   - Slightly negative residuals (-0.48) indicate recent underperformance relative to the model's expectations.

3. **Balanced Fundamental Picture**:

   - Strong positives include NVIDIA's dominant AI market position, inclusion in the "Magnificent 7," and potential benefits from the incoming Trump administration's "Project Stargate."

   - Significant concerns include valuation worries, competition from China's DeepSeek AI, and signs of speculation rather than fundamental-driven price action.

4. **Risk vs. Reward Assessment**:

   - While the positive drift suggests upside potential, the elevated tail risk metrics ($EVT_{VaR}$, $EVT_{CVaR}$) indicate substantial downside risk, particularly given the "AI bubble" concerns mentioned in the analyst thoughts.

   - The neutral RSI and current market uncertainty suggest waiting for more clarity before taking a directional position.

Given these factors, neither buying nor selling appears clearly advantageous at this moment. The positive drift and NVIDIA's strong market position argue against selling, while the elevated risk metrics and valuation concerns argue against buying. A HOLD recommendation allows for monitoring developments around competitive threats, valuation trends, and "Project Stargate" specifics before committing to a directional position.

Decision: **HOLD**

# B  Simudyne Horizon Synthetic News Context

## Summaries of Historic News Events

Example of historical news headlines:

> - (*2022-12-20*) The S&P 500 edged up marginally by 0.10% to 3,821.62 on Tuesday, managing to stay in positive territory despite broader global market weakness. Technology stocks showed some resilience, while investors appeared to be cautiously positioning ahead of Thursday's PCE inflation data and Friday's shortened pre-holiday trading session. The slight gain comes as the index continues to navigate concerns about the Federal Reserve's hawkish stance on interest rates heading into 2023.
>
> - (*2023-02-02*) The S&P 500 surged 1.44% to 4,179.76, extending its rally following the Federal Reserve's widely anticipated 25 basis point rate hike on Wednesday. Technology stocks led the gains as investors interpreted Fed Chair Powell's comments about "disinflation" as a signal that the central bank might be nearing the end of its tightening cycle, despite his caution that more increases would still be needed.
>
> - (*2023-03-09*) The S&P 500 tumbled 1.90% to 3918.32 on Thursday, marking its worst session in over two weeks as investors grew increasingly concerned about Federal Reserve Chair Jerome Powell's recent hawkish testimony to Congress. Technology and financial stocks led the decline after SVB Financial Group's announcement of a capital raise sparked broader concerns about bank balance sheets and liquidity in a rising rate environment. The labor market data released today showing higher-than-expected jobless claims further fueled investor anxiety about the Fed's aggressive tightening path.
>
> - (*2023-04-17*) The S&P 500 gained 0.32% on Monday, closing at 4,151.32 as investors showed renewed confidence ahead of key earnings reports from major tech companies this week. Financial and healthcare sectors led the advance, with market participants seemingly shrugging off recent banking sector concerns and focusing on the potential for a less aggressive Fed rate hike path following last week's cooler-than-expected inflation data.
>
> - (*2023-05-30*) The S&P 500 edged slightly higher on Tuesday, gaining a modest 0.05% to close at 4,205.52 as investors balanced ongoing debt ceiling negotiations with positive consumer confidence data. Technology stocks provided support to the index, with AI-related companies continuing their upward momentum following NVIDIA's recent rally. Market participants remained cautious ahead of upcoming inflation data and the Fed's next policy decision in June.

## Summaries of Synthetic News Events

Example of synthetic news headlines generated by `Simudyne Horizon`, which identifies news stories that are consistent with synthetic price paths. Further details are given in Section 2.2.3, and can be found at `the Simudyne Horizon` product page..

- (*2022-12-20*) The S&P 500 inched up 0.10% to 3,821.62 on Tuesday, showing minimal movement as investors remained cautious following the Federal Reserve's hawkish stance last week and amid ongoing concerns about the FTX cryptocurrency collapse's broader financial impact. Energy stocks provided modest support following reports of dwindling U.S. strategic petroleum reserves, while technology shares continued to face pressure as companies reassess spending plans in light of Amazon's recent massive layoff announcement. Market participants appeared hesitant to make significant moves heading into the holiday-shortened trading week, with many awaiting tomorrow's consumer confidence data for further insights into recession risks.

- (*2023-02-02*) The S&P 500 plunged 1.62% to 3,796.29 on Thursday as markets reeled from the widespread effects of a major cyber attack targeting global financial payment systems, now in its third day. Financial stocks led the decline as major banks and payment processors continued to struggle with service disruptions, while tech companies involved in cybersecurity saw mixed performance amid questions about their ability to contain the breach. Investor sentiment deteriorated further following testimony from Treasury officials suggesting the attack could have longer-term implications for financial stability than initially anticipated.

- (*2023-03-09*) The S&P 500 managed a modest gain of 0.51% to close at 3,800.86 on Thursday, showing resilience despite ongoing concerns about the banking sector following Credit Suisse's collapse earlier this week. Technology stocks led the rebound as investors cautiously returned to growth sectors after the Supreme Court's AI licensing decision provided some regulatory clarity, though financial stocks remained under pressure as market participants continued to assess potential systemic risks from the banking crisis. The index's performance suggests investors are beginning to look past immediate banking concerns while remaining cautious about the broader economic implications of recent financial sector turmoil.

- (*2023-04-17*) The S&P 500 tumbled 0.96% to 3,934.62 on Monday as investors reacted to the Federal Reserve's unexpected emergency meeting announcement scheduled for tomorrow, fueling speculation about potential interest rate actions outside the regular schedule. Banking stocks led the decline following reports that additional regional banks are facing liquidity pressures despite the recent Credit Suisse rescue, while technology shares were further pressured by ongoing uncertainty surrounding the Supreme Court's recent AI licensing requirements that continue to create implementation challenges for major tech firms.

- (*2023-05-30*) The S&P 500 edged up 0.25% to 4,226.27 on Tuesday, extending its cautious recovery as investors digested the recently announced bipartisan corporate tax reform proposal that has boosted market sentiment. Energy stocks led gains following rumors of verification data supporting last week's major North American oil discovery announcement, while technology shares showed mixed performance as the market continues to process the implications of the Supreme Court's AI licensing requirements. Healthcare stocks remained under pressure as analysts revised growth forecasts downward in the wake of last year's blocked pharmaceutical merger.