# From Deep Learning to LLMs: A survey of AI in Quantitative Investment

BOKAI CAO, The Hong Kong University of Science and Technology (Guangzhou), China and IDEA Research, International Digital Economy Academy, China

SAIZHUO WANG, The Hong Kong University of Science and Technology, Hong Kong and IDEA Research, International Digital Economy Academy, China

XINYI LIN, XIAOJUN WU, and HAOHAN ZHANG, The Hong Kong University of Science and Technology (Guangzhou), China and IDEA Research, International Digital Economy Academy, China

LIONEL M. NI, The Hong Kong University of Science and Technology (Guangzhou), China

JIAN GUO*, IDEA Research, International Digital Economy Academy, China

Quantitative investment (quant) is an emerging, technology-driven approach in asset management, increasingly shaped by advancements in artificial intelligence. Recent advances in deep learning and large language models (LLMs) for quant finance have improved predictive modeling and enabled agent-based automation, suggesting a potential paradigm shift in this field. In this survey, taking alpha strategy as a representative example, we explore how AI contributes to the quantitative investment pipeline. We first examine the early stage of quant research, centered on human-crafted features and traditional statistical models with an established alpha pipeline. We then discuss the rise of deep learning, which enabled scalable modeling across the entire pipeline from data processing to order execution. Building on this, we highlight the emerging role of LLMs in extending AI beyond prediction, empowering autonomous agents to process unstructured data, generate alphas, and support self-iterative workflows.

## 1 INTRODUCTION

Asset management is a crucial and expanding segment of the financial industry, with **Quantitative Investment (Quant)** emerging as a key approach within it. Quantitative investment strategies

---

*Corresponding Author.

Authors' addresses: Bokai Cao, mabkcao@connect.hkust-gz.edu.cn, The Hong Kong University of Science and Technology (Guangzhou), China and IDEA Research, International Digital Economy Academy, China; Saizhuo Wang, swangeh@ connect.ust.hk, The Hong Kong University of Science and Technology, Hong Kong and IDEA Research, International Digital Economy Academy, China; Xinyi Lin, xlin652@connect.hkust-gz.edu.cn; Xiaojun Wu, xwu647@connect.hkust-gz.edu.cn; Haohan Zhang, hzhang760@connect.hkust-gz.edu.cn, The Hong Kong University of Science and Technology (Guangzhou), China and IDEA Research, International Digital Economy Academy, China; Lionel M. Ni, The Hong Kong University of Science and Technology (Guangzhou), China, ni@ust.hk; Jian Guo, IDEA Research, International Digital Economy Academy, China, guojian@idea.edu.cn.

**111**

leverage statistical analysis, optimization techniques, and increasingly, AI algorithms to identify and exploit market inefficiencies. Benefiting from the exponential growth in data availability, computational power, and technological innovations, these approaches significantly improve investment decision-making and provide a competitive edge in the financial market.

Among various quantitative investment approaches, **alpha strategy** has received considerable attention for its strong capacity to capture market inefficiencies and its natural alignment with AI-driven predictive methods. The pursuit of 'alpha' refers to predicting individual asset's excess returns over the market's overall performance, such as a stock index, and is the central focus of portfolio managers. The development of alpha strategies typically includes four steps: data processing, model prediction, portfolio optimization, and order execution (as introduced in subsection 2.2). These four sub-tasks, though distinct, are closely interconnected, all working towards the common goal of maximizing excess returns while controlling risks. Compared to other quantitative investment strategies, such as high-frequency trading or arbitrage, alpha strategies have been shown to have great capacity and effectiveness by exploiting market mispricings. As a result, alpha strategies receive the highest attention, research focus, and market share of researchers and investors, representing the core technology in quantitative investment. In this survey, we take the alpha strategy as a representative example of quantitative investment and center our discussion on how AI plays a role in this field.

In recent years, the application of **deep learning (DL)** techniques in alpha strategies has shown promising results, demonstrating the ability to identify complex patterns and relationships in financial data that are difficult to detect using traditional quantitative methods. Meanwhile, **large language models (LLMs)**, such as GPT-series [4], BERT [38], and their financial variants, have shown remarkable power in understanding contextual data, generating accurate interpretations, and reasoning like human analysts. Therefore, their application in finance, particularly in quantitative investment, has sparked endless possibilities.

This paper focuses on the evolution, application, and respective advantages of DL and LLMs in quantitative investment, with a specific emphasis on alpha strategies, providing a comprehensive review of the existing literature, and discussing the potential, challenges, and limitations of how LLMs could enhance DL-based approaches.

## 1.1 Evolution of Alpha Strategy Investment

The evolution of the alpha strategy could been characterized by a three-stage progression from manual labeling of trading signals to the use of deep learning models and ultimately to an era of agent interaction and decision-making between LLM agents (Figure 1). In the early stages, the focus was on traditional statistical modeling of market patterns, relying on the expertise of individual researchers to identify profitable trading signals and develop corresponding models. However, this approach had limitations due to the complexity of financial markets and the difficulty of capturing all relevant factors in a model. It still relies heavily on the skill and experience of human researchers to evaluate and execute trading strategies.

As the field has matured, the application of deep learning has opened up new possibilities for quantitative investment. Particularly in the context of alpha research, deep learning has been shown to be effective in identifying inherent patterns. For example, deep learning models have been used to analyze factors such as spatial interconnectedness [142], long-term temporal dependence [178], and news sentiment [68] to predict the price movements and manage positions. While the use of deep learning holds great promise, there are also challenges that need to be addressed. One major challenge is the risk of overfitting, which can lead to poor performance when the model is applied to new data. Another challenge lies in improving interpretability and accuracy to further understand, reason, and interact with vast volumes of multimodal data. Despite these challenges,
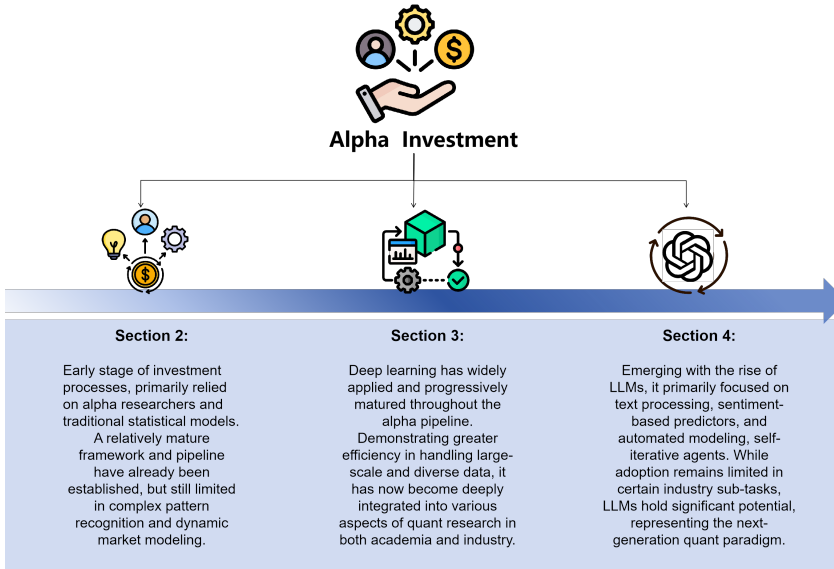
Fig. 1. The evolutionary process of Alpha investment across different stages.

the application of deep learning in quantitative investment is expected to continue to grow in both academia and industry, as investors seek to gain an edge in a competitive market.

More recently, LLMs have emerged as a powerful tool in quantitative investment, characterized by rapid development and enormous application potential, attracting significant attention. They excel in understanding and processing multimodal data and possess the potential to autonomously handle complex tasks of reception, comprehension, and inference over large-scale datasets. In current mainstream research, LLMs primarily serve two roles in alpha strategy: as predictors (subsection 4.1) and as agents (subsection 4.2), handling various tasks. In both cases, they contribute high-level insights built upon deep learning frameworks and have the potential to further evolve AI-powered deep learning investment methods into the AI-automated stage. However, the practical deployment of LLMs is still in its early stages. We will highlight their current limitations and discuss potential future directions for their development (subsection 4.3).

## 1.2 Motivation and Contribution of this Survey

The use of deep learning and LLMs in alpha strategies of quantitative investment has recently seen a surge of interest, with many studies focusing on the application in various aspects of the alpha research pipeline. However, most of these studies are relatively isolated in specific tasks or disciplines, and there is a lack of a unified view of the whole landscape of quantitative investment, particularly in the context of alpha strategy. This survey will also systematically summarize the evolution of alpha research from the perspective of phased algorithmic evolution. Additionally, quant is a field where research and practice are highly interconnected, while existing survey papers are limited to filling technical gaps between practical opportunities and research theories related to the combination of LLMs and DL-based alpha models. There's a lack of a comprehensive framework and forward-looking perspective on the future of deep learning and LLMs research workflows from a real-world viewpoint.

To address these issues, this survey paper aims to provide readers with a more integrated and comprehensive view of the alpha strategy. We intend to achieve this by surveying relevant works
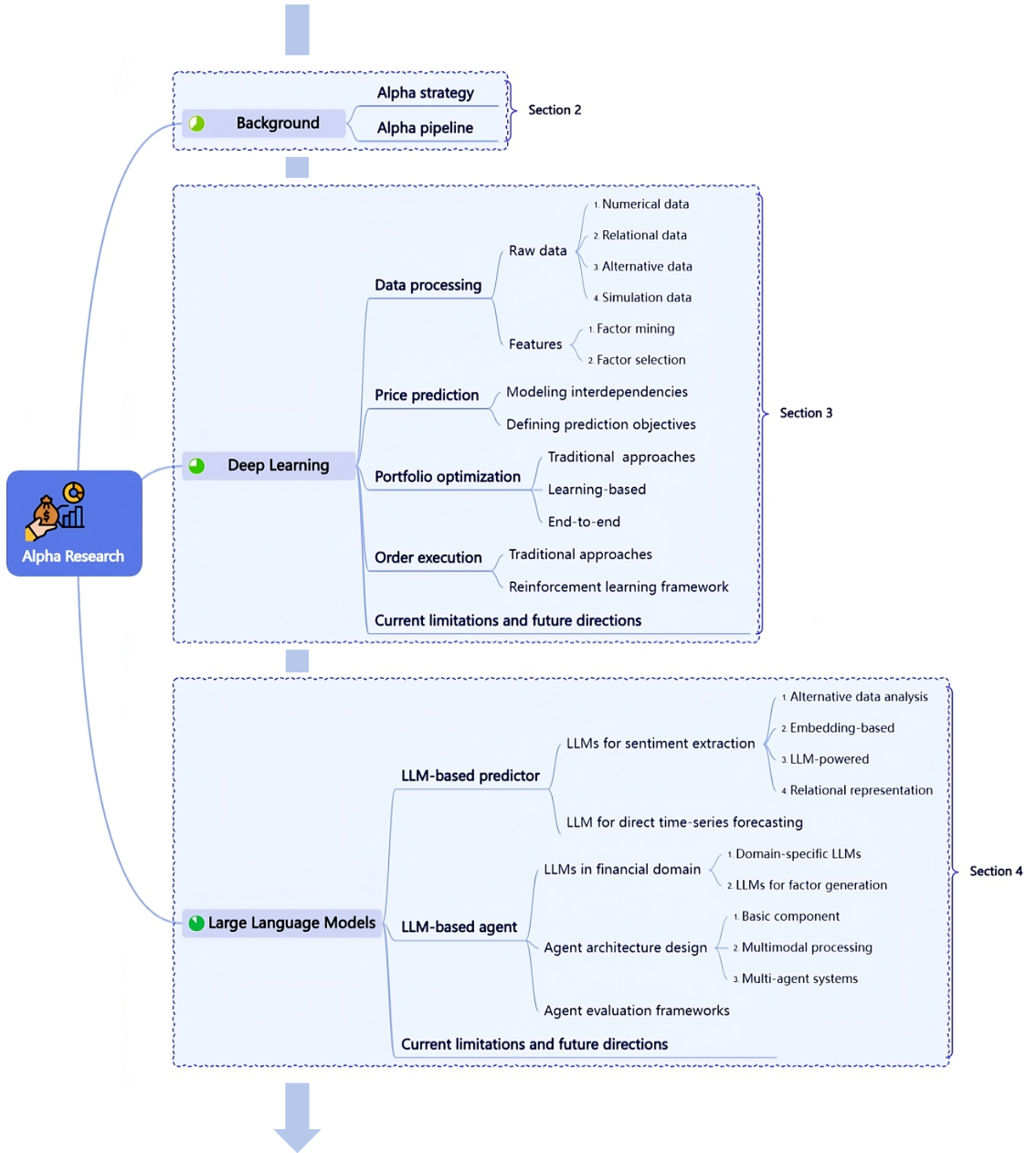
Fig. 2. The overall framework of this paper.

that cover all types of deep learning and LLMs tasks in the whole alpha pipeline, offering a holistic and interconnected view of the field. Moreover, our survey paper seeks to provide a broader research perspective by starting from real-world applications, highlighting the practical issues and challenges faced by investors to reveal the generic research problems and promote future studie. The overall framework of this paper is shown as Figure 2.

Key contributions of this paper include:

- It provides a comprehensive survey of the existing research on the use of deep learning and large language models in alpha strategies, connecting different works in a concrete research pipeline. This survey is the first of its kind to cover the topic comprehensively, offering a holistic view of the field.
- It Introduces the domain from an interdisciplinary standpoint, emphasizing practical applications to derive key research questions for quantitative investigations. Additionally, it discusses the most challenging problems from a practical perspective and offers insights into potential future research directions.
- It systematically compares the technical approaches, strengths, and weaknesses across the three stages of quantitative investment: traditional statistical models, DL-based methods, and LLM-based approaches. Building on the iterative development, it identifies key gaps and leads alpha strategies to the next stage.

## 2 BACKGROUND

In this section, we provide a brief overview of alpha strategy and its role in quantitative investment. Specifically, we begin by introducing alpha strategies and then discuss the alpha pipeline in quant, which serves as the framework for developing and implementing the investment process.

### 2.1 Alpha Strategy

Alpha strategies are investment approaches that focus on identifying opportunities yielding returns that exceed the market benchmark by exploiting inefficiencies or mispricings. These strategies typically consist of two components: the alpha side and the hedging side.

The alpha side of an alpha strategy focuses on generating excess returns by identifying profitable investment opportunities. To achieve this, it aims to make accurate predictions about the trends of individual instruments, sectors, or the overall market. These predictions are then used to allocate capital to different investment instruments to maximize returns. The alpha side faces several common challenges. For example, predicting asset prices involves handling high-dimensional, noisy data with complex patterns, while portfolio allocation requires optimizing the trade-offs between risk and return. Deep learning techniques have been shown to be effective in addressing these challenges, which is why recent works have applied deep learning to alpha strategies.

The hedging side of an alpha strategy focuses on managing the risks associated with market movements. Its goal is to mitigate market risks, ensuring that the excess returns generated by the alpha side are not eroded. To achieve this, the hedging side typically employs hedging portfolios, such as derivatives like stock index futures or options, while aiming to minimize hedging costs. The hedging side also faces several common challenges. For example, selecting the appropriate hedging portfolio requires understanding the relationship between the portfolio and the underlying investments. Additionally, hedging portfolios may incur transaction costs, which can affect returns.

While both the alpha side and the hedging side of an alpha strategy aim to generate returns while controlling risks, they face different challenges and require distinct solutions. In this paper, we primarily focus on the alpha side of alpha strategies and discuss recent works that have employed deep learning techniques and large language models to address the challenges faced by the alpha
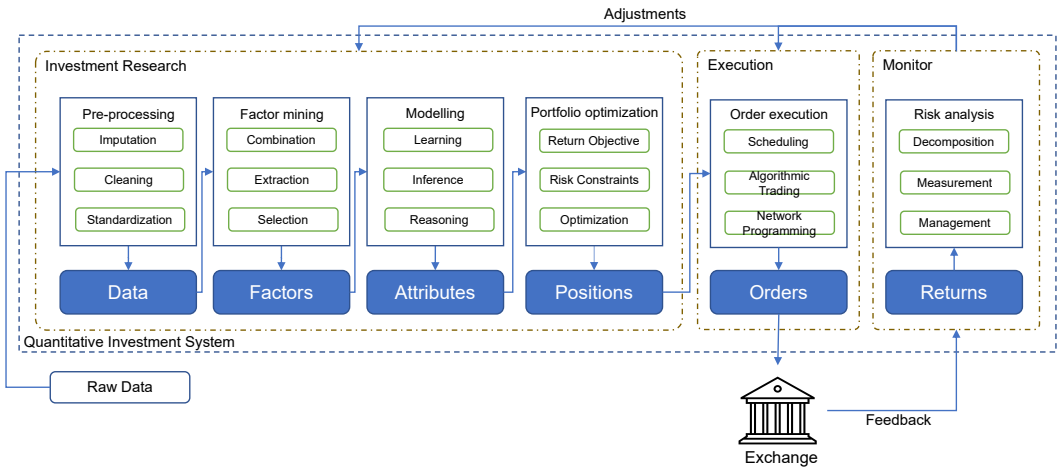
Fig. 3. A typical pipeline of quantitative investment.

side. However, it is important to note that the hedging side is also an essential component of alpha strategies, and interested readers can refer to relevant literature for further study.

## 2.2 Alpha Pipeline

In its essence, an alpha strategy aims to identify which instruments to trade, and actually trade them to earn profits. In its most simple form, this idea can be practiced in a one-step way. The investor can just look at the market, find some assets he wants to trade, and put the corresponding orders on the market.

However, things become much more complex when the capacity of the strategy grows in terms of both the amount of money and the number of assets involved in the investment. In this case, trading decisions cannot be carried out just with simple human operations. Instead, the whole strategy needs to be standardized into several sub-tasks, where each task is well-formulated and has a well-developed toolbox to deal with it. As shown in Figure 3, an alpha strategy can be decomposed into a pipeline consisting of several sub-tasks, which we will illustrate in the following.

- **Data processing**: The whole pipeline starts from analyzing the data from the market, and the data come in various forms. To perform modern data mining techniques on these data, the data must be first cleaned and standardized into unified forms. The pre-processing step is hence involved to do the cleaning, standardization, and imputation of data. Meanwhile, pre-processed data can developed into features to better integrate market information and serve as input for the next stage of the model.
- **Model prediction**: Although investments can be done with various motivations, the most important one should be the expected future performance of the asset price. Therefore, to form an investment decision, we need to first predict the future of the assets of interest, and then make decisions accordingly. Price prediction aims to make predictions about assets, such as their future price change, volatility, etc. This prediction task is what deep learning is good at, so a bunch of deep learning techniques have been studied in making better and more accurate predictions.
- **Portfolio optimization**: Model predictions themselves cannot be directly used as investment decisions. They are utilized in the portfolio optimization stage to generate investment decisions. Essentially, the portfolio optimization step takes the various kinds of market predictions (e.g.

Table 1. Categorization of financial data.

| Modality | | | Features | Example | References |
|---|---|---|---|---|---|
| Numerical | Quote data | Regular interval | Quotes that are generated at regular time intervals | 1-minute candlestick chart | [131, 149] |
| | | Irregular interval | Quotes that are generated at irregular time intervals | Tick-level[1] order book data | [16] |
| | Fundamental data | | Fundamental data such as revenues and profits. | Financial statements | [134] |
| Relational | Pairwise edges | | The relation between a pair of entities | Knowledge graph | [46, 75, 159, 160] |
| | Hyperedges | | The relation involving a set of entities | Sector categorization | [127] |
| Alternative | Text | | Information expressed in natural language. | Social media posts | [59, 87, 128, 162] |
| | Images | | Images that are related to the traded asset | Satellite images | [101, 118] |
| | Other modalities | | Anything | WiFi traffics, cell phone signals | [64] |
| Simulation | Time series | | Synthetic quote data | Simulated orders | [50, 84, 183, 184] |
| | Tabular | | Database structured in a tabular form | Simulated financial statements | [47, 126, 133] |

price change, volatility) as input, and outputs the corresponding investment decisions such as the amount of money allocated for each asset in the next holding period (also known as *positions*). The core problem in this stage is to find an optimal portfolio optimization that maximizes the objective defined by investors such as maximizing the expected risk-adjusted returns, subject to the constraints defined simultaneously, such as the maximum volatility constraint or the diversity constraint defined for risk control. Traditionally, this problem is formulated as an optimization problem and the corresponding optimization techniques can be applied to solve this problem.

• **Order execution**: The portfolio allocations need to be implemented by actually putting orders on the market and making the deal. And this order execution step is by no means an easy task for alpha strategies that involve a large amount of money. This is because orders placed on the market will inevitably bring fluctuations on the market, driving the price to the opposite direction that is beneficial for the investor. And when the trade volume becomes large, this effect is also magnified and might introduce great loss to the strategy. Hence, the order execution stage is involved to minimize such loss by splitting big orders into smaller ones and deciding the appropriate time to execute them. Traditionally, this problem has been formulated as an optimal control problem where the dynamics of the limit order book are extensively studied.

The investment pipeline is not an open-loop system. Instead, it needs feedback from the market to adjust its behaviors. The monitor module is therefore involved in the entire pipeline to conduct risk analysis, collect feedback signals, and make corresponding adjustments to the previous processes. For example, the risk management module calculates the risk exposure of the current portfolio to different sectors. Then, from the perspective of the overall process, it makes balancing or neutralization adjustments by injecting constraints or adjusting the objective functions in the portfolio optimization module.

## 3 DEEP LEARNING IN ALPHA PIPELINE

Deep learning has been widely applied throughout the whole alpha pipeline. In this section, we systematically analyze how deep learning is utilized to enhance traditional alpha research at each sub-task of the previously discussed pipeline.
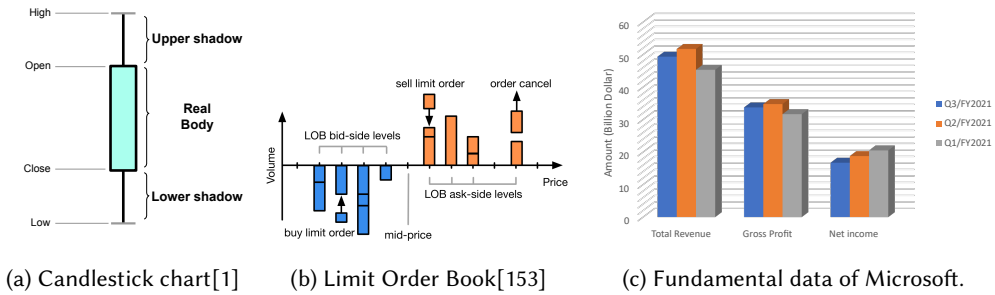
(a) Candlestick chart[1]    (b) Limit Order Book[153]    (c) Fundamental data of Microsoft.

Fig. 4.  Examples of numerical data.

## 3.1  Data Processing

As stated before, the data used for alpha research need to be first pre-processed into unified format, and then used for prediction. Here we divide data before and after pre-processing into two types, namely raw data, and features.

*3.1.1  Raw Data.* Financial markets constantly generate heterogeneous data with various modalities. Based on their source and modality, we categorize financial data used in quant strategies as numerical data, relational data, alternative data and simulation data. A summary and a comparison of different types of data are presented in Table 1 and we will elaborate on each of them in the following.

*Numerical data.* Numerical data are the most widespread type of data in the financial world, it can be categorized as quote data that indicate the movement of asset prices of arbitrary frequency, and fundamental data that reflect the operations of the underlying economic entities of securities. Quote data typically includes candlestick chart and limit order book, as illustrated in Figure 4. Candlestick chart is a way to illustrate the price movement of a financial instrument, which consists prices at four different information dimension in an time interval, namely the open, close, high, and low prices. A limit order book refers to an electronic list of buy and sell limit orders organized by price levels. Fundamental data is widely used by analysts to determine the intrinsic value of a financial instrument. Important sources of fundamental data are balance sheets, income statements, and cash flow statements from financial statements.

Financial markets are intrinsically dynamic, numerical data are often treated as time series. Under such circumstances, a very important property is their frequency, i.e., the interval at which new data points are generated. Generally speaking, fundamental data have the lowest frequency, changing every few months. In contrast, quote data have frequencies at multiple levels, ranging from the lower ones such as week- or day-level to the higher ones such as minute- or second-level frequencies. More extremely, the quote data with the maximum temporal resolution, namely tick-level data, records the history of each order and transaction sequentially.

*Relational data.* While numerical data describes individual financial entities, the relationships between them can also impact market trends. We refer to such data as *relational data*, which describes the ubiquitous relationships between two or more financial entities. Formally, relational data are usually represented as a graph $G = V \times E$, where $V$ is the set of nodes and $E$ is the set of edges. The edges can be either *pairwise edges* or *hyperedges* [147], based on the number of entities involved in the relations. In practice, relational data have rich semantics, ranging from the business relationships between companies, to various events involving financial entities, to the correlations

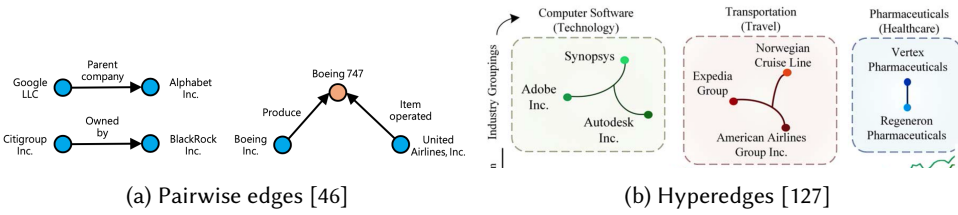(a) Pairwise edges [46]          (b) Hyperedges [127]

Fig. 5. Relational data examples.

and causal relationships between entities from a statistical perspective. An illustration of pairwise edges and hyperedges is present in Figure 5.

(1) Pairwise edges: As the most common type of relations, pairwise edges describe the relations between a pair of entities. A pairwise edge can be represented as a triple $(v_1, r, v_2)$ where $v_1, v_2$ are entities and $r$ is the edge type or edge weight. Meanwhile, it can be further categorized into static and dynamic depending on whether the underlying relationship changes over time. Static edges usually include the stable relationships such as upstream/downstream partnerships. Dynamic edges usually describes the event-based relationships that happens at arbitrary time points.

(2) Hyperedges: Apart from pairwise edges, financial entities are also involved in set-based relations, which are expressed as hyperedges. Assets like stocks are often categorized based on their sectors, such as technology, real estate and healthcare. They are also segmented based on related concepts [94]. For example, the concept "Metaverse" is related to a set of stocks such as the Meta Platform, Nvidia, Unity, and Roblox.[2] Formally, a hyperedge $\mathcal{R} = \{n_i\} \subseteq V$ is a subset of $V$ involving a few entities. Graphs containing hyperedges are termed *hypergraphs*. As a natural extension of graph neural networks, hypergraph neural networks are also being actively studied both theoretically and practically.

*Alternative data.* Alternative data refer to the multimodal data, such as text, image, and speech that convey predictive information. Quant models can refine their investment decisions by tapping into alternative data, which offers unconventional insights into diverse perspectives about the financial market. For example, the news about Elon Musk owning 9.2% of Twitter makes the price of Twitter close up 27%.[3] The success of deep learning in domains such as computer vision and natural language processing brings us tools to make decisions based on such comprehensive, multi-modal information. We can apply deep learning to extract useful events from news [123], predict the sentiment of a post [95], *mining the lead-lag relationship from knowledge graph [105]* , and count the number of customers visiting Costco in a day [31].

*Simulation data.* To improve quantitative models, high-quality data are essential. However, acquiring large-scale real-world financial data is challenging due to limitations in limited availability, privacy concerns, and high costs [70]. In this context, simulation data offer a viable solution. Synthetic data enable better incremental training, robustness assessment, and risk testing for alpha model development. The main generation methods fall into four categories:

(1) Rule-based:
Early market simulation studies heavily relied on rule-based approaches, where predefined rules dictated trading behavior under specific price conditions, often assuming trend following or mean reversion [116, 119], or simplistic resample methods. While these models offer interpretability and

---

[2]https://bit.ly/3OP5QRu

[3]https://bloom.bg/3bzOABA

ease of implementation, their rigidity limits adaptability to real-world market dynamics, which involve non-linear relationships and evolving dependencies [70].

(2) Time Series-based: Time series-based methods use historical data to simulate market dynamics. The advancement of deep learning enables the capture of complex financial patterns, including Variational Autoencoders (VAEs) [76, 122], Generative Adversarial Networks (GANs) [50, 52, 183] and diffusion models [60]. These models facilitate market trajectory prediction [29, 30] and the simulation of market fluctuations and risks [32, 139, 150, 154]. While effective in preserving statistical properties, they still struggle to model causal interactions between market participants.

(3) Order Flow-based: With improvements in market microstructure research and computational efficiency, order flow-based methods are emerging. These models simulate order arrivals, executions, cancellations, and limit order book dynamics to capture price formation and volatility [27, 28, 85]. They aim to uncover market dynamics by modeling order interactions, with examples like DeepLOB [184], MarS [84], and DiGA [61]. Challenges include modeling complexity, high computational costs, and difficulty capturing long-term economic cycles beyond short-term fluctuations.

(4) Multi-Agent-based: Recent simulation methods also focused on multiple agents, modeling traders, investors, and dealers with distinct goals and decision strategies [18, 102, 173]. These agents interact with each other and the market, producing emergent phenomena that mirror real-world markets. While capturing market complexity, these approaches struggle with realism-efficiency trade-offs and parameter calibration for long-term accuracy [146].

*3.1.2 Features.* Financial data are intrinsically noisy and large-scale [13, 104], making it difficult to extract meaningful information directly. Hence, we use the attributes derived from the original data, namely factors [2, 43, 124], to describe the asset from different financial aspects such as value, size, momentum [4] [22], reversal [5], volatility [6], etc. Mining factors integrate information, forming the foundation for subsequent quantitative models [136]. From the perspective of deep learning, the factor mining step corresponds to feature engineering [3] for financial data, the typical workflow consists of three procedures including feature construction, feature extraction, and feature selection [93].

Feature construction aims to enhance model performance by creating more informative representations, known as factors, through the strategic combination and manipulation of raw data. Financial factors can be categorized into two types: 1) Symbolic factors, expressed through symbolic equations or rules [69], and 2) Machine learning factors, derived from machine learning models [176].

Traditionally, symbolic factors have been explored by human researchers, relying heavily on expert knowledge and intensive labor. To broaden the scope of exploration and expedite the factor mining process, researchers are increasingly turning to algorithmic factor mining. This approach, essentially a task of symbol regression, can be addressed with regression techniques such as genetic programming algorithms [14, 121, 179] and neural symbol regression [34, 167].

While symbolic factors offer higher readability, their effectiveness is often limited by the predefined scope of operand and operator space, convergence efficiency and factor diversity. By contrast, machine learning factors boast a more flexible representation, owing to their high-dimensional parameter space. This flexibility allows for a more nuanced and adaptive approach to factor construction in financial modeling. Researchers typically employ feature extraction models, such as encoder-decoder architectures [42], to align with predetermined learning objectives, like predicting

---

[4]https://bit.ly/2HZHofh
[5]https://bit.ly/3PxLkFG
[6]https://bit.ly/2QoH8Oq

(a) Temporal Patterns                (b) Spatial Patterns                (c) SpatioTemporal Patterns.
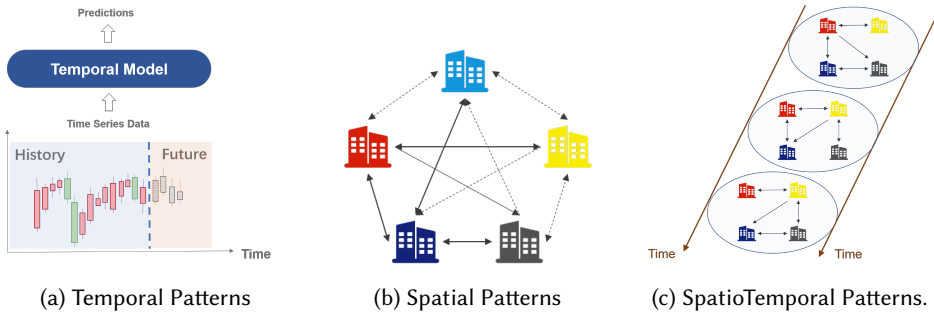
Fig. 6. Modeling Data Interdependencies.

future returns or volatility. In this process, either the latent representation or the prediction output of the model can be utilized as a signal factor [158, 161]. The effectiveness of machine learning factors stems from their robust expressive capabilities, which are attributed to high-dimensional parameters and non-linear functions, coupled with faster search efficiency directed by target gradient optimization. However, the primary challenge with these factors lies in their limited interpretability and readability, which poses difficulties in risk management.

The goal of feature selection is to select a subset of features for dimension reduction, avoiding overfitting and improving model performance. For quantitative investment, this step is usually intended for selecting the best-performing factors individually. Hence, one commonly used method for feature selection is the filtering method [37, 54, 78], which scores each feature according to certain criteria and the best-performing ones are retained. As for the criteria, the correlation between the factor and the actual return [174] is often used to measure its predictive capability. The correlations between selected factors should be minimized to eliminate redundant information.

It should be noted that factors are not always necessary. Some approaches directly model from raw data inputs in an end-to-end fashion. With the advancements in deep learning techniques, these methods are increasingly favored by researchers, we will will elaborate on this in subsection 3.2.2.

## 3.2 Model Prediction

In the field of prediction, researchers have extensively employed various deep learning methods, many of which have demonstrated practical effectiveness in real-world markets. The essence of deep learning modeling lies in two core aspects: the model architecture and the optimization objective. The model architecture is selected based on the inherent relationships between the targets and the available data, and it must possess the expressive power to capture these complex interdependencies. Meanwhile, the optimization objective determines how the quality of the model's output is evaluated and guides the training process toward effective predictions. For the task of price prediction, it is crucial to address two fundamental questions:

(1) How is future price information embedded in historical data, and how can we effectively model these relationships?

(2) How to define the optimization objective so that the prediction signals best support subsequent trading strategies?

*3.2.1 Modeling Data Interdependencies.* Financial data inherently exhibits temporal and spatial correlations. Temporal correlations (Fig 6a) reflect how data points are related over time, such

as trends, momentum, or reversals in market dynamics. Spatial correlations are about the inter-connectedness between entities, such as stock interdependencies, market sector influences, and upstream-downstream relationships in industry supply chains. Asset prices are driven by a complex interplay of these temporal and spatial factors. Accurately modeling these dynamics can lead to improved forecasting, as it allows predictive models to harness both historical trends and cross-sectional relationships among market entities. To achieve this, researchers have developed diverse models that capture these complex interactions. Depending on the nature of the relationships they aim to capture, prediction model architectures are typically categorized as temporal models, spatial models, or spatiotemporal models.

*To capture temporal patterns.* Modeling temporal patterns is essential because it enables us to extract valuable information from historical data, detect trends, and forecast future price movements based on time-evolving signals. Therefore, to capture these time-dependent dynamics, researchers construct time-series inputs using historical price and volume data from assets and apply temporal modeling techniques that aggregate information across successive time steps. These models typically rely on the assumption of time translation invariance, meaning that the aggregation rules and parameters remain consistent across different time intervals. Typical examples of temporal blocks include convolutional neural networks [24, 57, 112, 132], RNNs [108, 177], and transformers [148], along with their respective variations. Moreover, many studies employ hybrid combinations of these blocks—for instance, CNN-LSTM[111] architectures—to capture temporal dependencies across different receptive fields, effectively modeling both local and sequential temporal correlations. Specialized models are also devised for irregular time intervals. For example, [185] uses fine-grained feature-level time span information to decay the effect of previous timesteps for making use of irregular time intervals. [66] partition the order sequence into discrete segments and perform temporal signal extraction on these segments to effectively model market microstructure.

*To capture spatial patterns.* Capturing spatial patterns is crucial because it allows models to leverage the relationships and dependencies among different assets and their corresponding sectors. By understanding how different entities interact, models can improve predictions by considering the influence of correlated market movements and sector-wide trends. To capture these spatial relationships, researchers adopt spatial modeling techniques that can be divided into two primary approaches: implicit and explicit methods. **Implicit methods** typically use self-attention mechanisms that evaluate the entire set of entities simultaneously without relying on predefined graph structures [46]. In contrast, **explicit methods** represent these relationships using sparse graph structures—either constructed directly [87, 159], or inferred via graph structure learning methods [159, 186]. Graph neural networks (GNNs) [77, 159] are then applied on the graphs to discern meaningful patterns. GNNs are particularly effective at handling complex graph structures, including those that involve hyperedges [127], by facilitating message passing across nodes. Most graph neural networks function by facilitating message passing on the graph. In addition, some researchers [36, 166] have explored the use of generative models, such as diffusion models, to generate dynamic asset graph structures that simulate the complex, time-varying relationships among different assets. Spatial methods typically assume spatial position invariance at different scales. For instance, self-attention maintains invariance to permutations in entities' positions, whereas GNNs ensure node and edge permutation invariance, allowing them to operate on various relational structures. Global and local methods in modeling have complementary strengths in their receptive fields. [74] using Graph Attention Networks (GATs) for large-scale financial graph structures combines graph modeling of prior relationships with attention mechanisms. This dual approach helps focus on crucial nodes, diminishing the influence of complex background noise and enhancing the signal-to-noise ratio.

*To capture spatiotemporal interactions.* Modeling spatiotemporal interactions is critical because it enables models to capture both the evolution of individual assets over time and the dynamic relationships between different market entities simultaneously. This comprehensive approach improves predictive accuracy by integrating insights from both dimensions. To capture these dual dynamics, researchers have adopted spatiotemporal modeling techniques that fuse spatial and temporal information. There are primarily two ways to achieve this integration: the decoupled and coupled approaches. The **decoupled approach** independently encodes spatial and temporal features; for example, a hypergraph encoder might first be applied to capture spatial relationships, followed by gated temporal convolution to model time-series dynamics [58]. MATCC[20] also designed a correlation module composed of multiple layers of attention and mixer junction submodules, each dedicated to modeling temporal correlations, inter-asset relationships, market trends, and other related factors. In contrast, the **coupled approach** integrates spatial and temporal information concurrently, thereby capturing the interactions between the two dimensions more directly and yielding improved representations of market behavior [127, 140].

3.2.2 *Defining Prediction Objectives.* Once a model architecture is established, it is equally important to define an appropriate optimization objective that guides the model toward effective and accurate predictions. In quantitative investment models, defining learning objectives is particularly challenging because key components of the investment workflow—such as portfolio optimization and order execution—are typically non-differentiable. To address this, prior research has generally classified training objectives into two main categories: intermediary targets and direct end-goal optimization. Intermediary targets primarily assess future price movements of stocks, subsequently utilizing these outputs as predictive signals for portfolio construction through optimization algorithms. The alternative approach focuses directly on optimizing the aggregate performance of the final investment portfolio, thereby ensuring alignment of the output with the ultimate investment goal. Each approach offers its own benefits and trade-offs for aligning model outputs with real-world trading objectives.

*Intermediary targets.* Opting for intermediary targets as the primary objective in model building offers simplicity by negating the need to consider complex elements like downstream risk control and the indifferential nature of order execution. However, this approach necessitates close coordination with downstream tasks, potentially leading to the accumulation of errors or a misalignment between the predicted targets and actual trading requirements. In modeling price movements, two main approaches are employed: individual trend analysis and relative ranking. The individual trend method focuses solely on the future price trajectory of a single asset, typically using the asset's future returns as the label. This approach can be modeled either as a regression or a classification task. Regression tasks [177], leveraging mean squared error as the loss function, provide more precise fitting values for expected future returns, thereby being more conducive to downstream tasks. Classification [127], on the other hand, segments return labels into multiple categories, like price increases or decreases, commonly using cross-entropy as the loss function. Given the inherently low signal-to-noise ratio in financial data, numerous studies have explored ways to enhance signal quality through techniques like data sampling and label denoising. For instance, the LARA framework proposed in [172] employs locality-aware attention to extract more informative samples from the data and uses RA-Labeling to adjust the labels of noisy samples on a per-trade basis during training, thereby improving the predictor's accuracy. While these methods smooth out data noise, it often results in a larger gap with subsequent optimization tasks, necessitating additional signal transformation. Relative ranking [46], in contrast, focuses on an asset's relative position within a cross-sectional framework. This method is particularly synergistic with certain portfolio strategies, such as long-short hedging strategies. Loss functions in relative ranking are classified into two

types: local ranking, exemplified by pair-wise ranking methods[41] which evaluate pairs of assets to determine superior future performance, and global ranking, which involves inputting a group of assets and optimizing for the correlation between predictions and actual future returns, or other metrics pertinent to learning-to-rank tasks.

*End-goal optimization.* End-goal optimization in portfolio modeling directly addresses final portfolio positions [96], focusing on complexities like inter-stock position control, risk management, and return stability from the start. This approach, while complex, effectively minimizes error accumulation common in multi-step models. Models employing end-goal optimization analyze a group of assets as a single sample, using spatio-temporal modeling approaches to account for both time-series trends and inter-asset relationships. Outputs typically include multi-period portfolio positions, optimized against performance metrics like return rates and the Sharpe ratio. Despite its complexity, this direct approach ensures alignment with portfolio management goals for comprehensive investment strategy. Despite the potential of this method, current research is limited by data availability and task complexity, leading to modest outcomes so far. However, with advancements in large-scale model technologies across various fields, this end-to-end modeling approach holds significant potential for future breakthroughs.

## 3.3 Portfolio Optimization

Portfolio optimization seeks optimal asset allocation to balance expected returns against volatility, transforming predictions of asset states into actual portfolio construction. This task, traditionally reliant on return and volatility predictions from statistical models, has been extensively researched within mathematical finance and operations research, as discussed in subsection 3.3.1. With the advent of deep learning, the field has witnessed a transformative shift towards data-driven approaches, leveraging the vast availability of financial data and computational advancements. Deep learning has been applied in two main ways: enhancing existing optimization components (subsection 3.3.2) and pioneering end-to-end methodologies for direct allocation generation (subsection 3.3.3). These advancements signify a new era in portfolio optimization, combining traditional insights with the capabilities of modern deep learning.

*3.3.1 Traditional Approaches.* Portfolio optimization began with Markowitz's Modern Portfolio Theory, which formulates optimal portfolio generation as a quadratic programming challenge. Subsequent methodologies, such as those based on Kelly's criterion, aim to maximize cumulative returns over multiple periods. This subsection divides the problem into two scenarios: single-period portfolio selection, focusing on MPT [80], and multi-period portfolio issues, concentrating on the latter [81].

*Single-period portfolio: Mean-variance Approach.* A single-period portfolio optimization problem can be simplified as generating a position for next holding period that can maximize the expected return while minimizing the potential risk. The expected return is usually represented as the mean of asset returns and risk is represented as the various of these returns, and such notation leads to its name: mean-variance approach, which is intended to balance return and risk.

The framework proposed by Markowitz is a very simple framework, and it cannot accommodate many practical considerations. Therefore, many follow-up works have been proposed to improve the original framework. Some important problems include the following.

(1) Adding regularization terms into the optimization objective to improve the robustness of the portfolio or reach new goals [21, 33]. For example, control of transaction costs can be realized by adding regularization terms to minimize the turnover rate. The position size can also be regularized by injecting sparsity regularization.

(2) Encourage diversity in portfolio allocation to reduce potential risk [129].
(3) Better estimate the covariance matrix. As risk measure, the covariance matrix of assets is usually estimated from historical data. However, such estimation can lead to a severe problem that the observations (e.g. the number of trading days) used for estimation are insufficient to generate reliable covariance. Hence, various methods have been proposed to address this issue [106, 115], such as many estimators proposed in statistics, and the factor model that leverages dimensionaility reduction to explain asset returns using a small number of factors.
(4) Better risk measures. Covariance has several limitations: it is hard to estimate, and it is not always an ideal risk measure that generalizes to every scenario. Hence, other risk measures have also been proposed, including value at risk (VaR), conditional value at risk (CVaR).

*Multi-period portfolio: Online learning and stochastic control.* In the multi-period setting, the focus has changed to maximizing the cumulative return across multiple holding periods. Instead of generating one portfolio vector that has been extensively polished, now we need to generate a series of portfolio allocations whose cumulative returns across multiple periods are maximized. Relevant techniques include:

(1) Online learning with heuristics: the portfolio allocation for current period can be computed via optimization. The optimization objective can be formulated based on some basic trading ideas, such as momentum (follow-the-winner) and mean-reversion effects (follow-the-loser)[81].
(2) Trend representation and pattern matching: naively following the basic ideas may not be effective in some cases. Hence more complicated trend representations are proposed to indicate some new trading ideas. Based on the predicted asset patterns and market distribution, new optimization objectives can be formulated.

On the other hand, if we can model the dynamics of asset price, then the multi-period portfolio optimization problem can be regarded as an optimal control problem, where the cost function to minimize is the negative cumulative return. However, in practice the exact dynamics are usually impossible to be accurately modeled, so uncertainty is introduced to allow for better flexibility, and this problem now becomes one of stochastic control.

*3.3.2 Learning-based Portfolio Optimization.* Given that market dynamics and patterns are stochastic and difficult to predict, deep learning can be incorporated into the framework presented above to enhance certain modules, namely improving the dynamics, the estimator, and the solver.

(1) [63] run Markowitz on specific returns (do neutralization on market risk factors via 'spectral extraction'), and predict the distribution of the specific returns to compute their expected mean and variance. After getting the mean and variance, a Markowitz can be conducted.
(2) [103] use neural networks for the prediction of returns and risks and apply traditional portfolio optimization (Markowitz/Omega) to generate the final positions.
(3) [141] studies continuous-time continuous-action and continuous-space portfolio optimization with finite horizon formalized as a stochastic control problem, and a Gaussian policy with time-decaying variance is derived. The Gaussian-based policy is proven to be better than adaptive control algorithms and DNN RL trained using DDPG. The asset price is modelled under GBM (Geometric Brownian Motion) and hence the dynamics are described by SDE. It is hence a stochastic control problem.

*3.3.3 End-to-end Portfolio Generation.* In addition to some works that use deep neural networks to combine parameter estimation and portfolio construction to achieve an "end-to-end" effect [8, 137], the reinforcement learning model has achieved better market performance in the field of end-to-end portfolio optimization itself. Since there is no universal, fixed and explicit label in portfolio optimization, it eliminates the need for nontrivial label construction, in order to achieve

a more flexible risk-return balance. Without the need for complex model and parameter tuning, reinforcement learning models are trained on direct trading end-to-end feedbacks. In this case, designing a reasonable objective function becomes important. We now discuss several methods to design these functions.

*Return-only approaches.* Considering only portfolio returns is the most intuitive way of modeling. [65] extract features from cross-sectional data to get the score of each asset, and then perform a normalization to get a position vector, the reward is then computed as the dot-product between this position vector and the return.

*Risk-adjusted returns.* [143] considers risks in investment and incorporates the Sharpe ratio as the reward function. The whole model is trained via policy gradient methods by propagating the Sharpe loss to different model parts.

*Transaction costs.* [182] considers the additional cost brought up by consecutive position adjustments, and added a regularization term computed as the l1 norm between the differences of the positions between 2 neighboring ticks.

*Diversity.* [110] considers different portfolio management styles may have different strengths in different markets. In this way a meta learning strategy is used to select from a number of different portfolio models trained with datasets generated by trading experts with different styles.

## 3.4 Order Execution

The order execution system implements the results from portfolio optimization, serving as a bridge between theoretical calculations and actual positions. It involves strategically placing and completing orders with the objective of minimizing total trading costs, considering factors including market conditions, asset liquidity and order impact. Research on order execution is usually conducted on high-frequency data, such as limit order book (LOB). The actions taken by participants in financial markets have become increasingly based on quantitative analysis and algorithms rather than any human decision making process [40], the prevailing methods could be delineated into traditional optimal control models and reinforcement learning models.

*3.4.1 Traditional Approaches.* Assuming that the dynamics of the limit order book, including both its inherent dynamics and the impact of market orders, can be analytically represented as the problem of executing with minimum cost. This optimal control framework can be divided into discrete and continuous models, both of which aim to find a dynamic trading strategy that executes transactions over a fixed period, with an optimal utility function reflecting a combination of cost and volatility.

*Discrete model.* The fundamental discrete trading model, proposed by the works of Bertimas & Lo [12] and Almgren-Chriss [7], assumes the price dynamics to be a discrete arithmetic Brownian Motion with fixed trading cost. They assumed that market impact is the only endogenous factor and the price volatility is exogenous. It is assumed to be the result of market forces that occur randomly and independently during trading, since market participants could detect and adjust their order placement.

Taking the revenue volatility of different trading strategies into account, there is a trade-off between impact and variance. The impact of risk on optimal execution could be constructed by solving the minimization problem of the expectation of shortfall for a given level of variance. For each level of risk aversion, there will be a uniquely determined optimal execution strategy expressed as the efficient frontier of optimal execution.

*Continuous model.* Continuous model further magnifies the time steps of the discrete model to infinity, solving the Almgren-Chriss problem with the dynamic asset price following specific distribution like Geometric Brownian Motion. The optimization of this stochastic control system can be transferred into a linear quadratic problem on a complete filtered probability space. [72] considers discrete trading under continuous GBM process, [48] provides numerical solutions under GBM using a mean-quadratic-variation objective function, [11] proves optimality under geometric price dynamics.

The above model operates under the assumption of a zero - drift random - walk price process. To achieve better price dynamics, it explores three approaches for scenarios where the underlying assumption fails. First, the price process might exhibit drift, indicating a potential directional view of traders. Second, it could display serial correlation. Third, a future regime shift due to events such as earnings announcements is anticipated. Based on these assumptions, some studies generalize initial conditions using more realistic price dynamics. Almgren [6] proposes a non-linear stochastic price impact model since the previous linearity assumption [7] deviate a lot from factual accuracy. [5] solves the optimal strategy numerically under specific assumption of market liquidity and volatility. More realistic representation of price dynamics including temporary mean-reverting [49] are also introduced to refine stochastic price impact to find the optimal strategies.

### 3.4.2 Reinforcement Learning Framework.
Reinforcement learning (RL) is effective in optimizing order execution due to its adaptability in dynamic, sequential decision-making environments. It enables end-to-end adaptation to changing market conditions by learning from the rewards associated with executed actions without relying on preliminary oversimplified model assumptions. Task components of RL are highly suited to the order execution process, corresponding relationship is as follows:

- State: Market variables of execution process, such as remaining time and current position.
- Action: A decision set of trading process, usually consists of submit, modify or cancel orders.
- Reward: The transaction cost and impact with a market order.

The reinforcement learning agent structure was first proposed by [107], which used Q-learning algorithms to train an optimized strategy for finding different price levels in given limit order book data. [56] adapted the solution from the Almgren-Chriss model to a dynamic strategy considering the specific current market microstructure using Q-learning. [89] proposed a Deep Q-Network (DQN) based RL algorithm for order splitting, and [109] applied a zero-ending inventory constraint to a double-DQN method.

In terms of policy-based algorithms, [35] was an early work that studied Proximal Policy Optimization (PPO) in order execution tasks, distinct from DQN. Numerical experiments were conducted in different environments with various dynamics. [44] proposed an actor-critic style policy distillation algorithm, and the model was trained on multiple assets to obtain a universal trading strategy. [55] used deep policy gradient methods to optimize strategies based on linear quadratic regulator problems. [90] presented a more end-to-end approach, using PPO with the limit order book and inventory as the state to directly output actions.

There are also model-based and multi-agent solutions. [91] established a multi-agent simulation system for order execution, where agents were trained and compared using various RL methods. An empirical study in [92] demonstrated the generalizability of DRL-based order execution agents, showing that different RL agents can transfer quite well.

The efficacy of different reinforcement learning algorithms has been validated across authentic datasets from multiple markets. The heightened focus and expeditious advancements of reinforcement learning aims to enhance order execution decisions across diverse financial markets,

particularly in scenarios where participants possess restricted information regarding the market dynamics.

## 3.5 Current Limitations and Future Directions

As the field of deep learning in alpha research continues to evolve, several emerging trends and areas of research promise to further revolutionize this domain. These advancements aim not only to enhance predictive accuracy and reliability of exeuction models but also to streamline the investment process, improve interpretability, and adapt to the dynamic nature of financial markets. Below, we outline four pivotal areas that encounter limitations but represent promising avenues for future exploration.

*Automated Machine Learning.* The development and tuning of predictive models for financial markets are both labor-intensive and complex, exacerbated by the markets' dynamic nature which necessitates frequent model updates. Automated Machine Learning (AutoML) presents a promising solution to streamline this process, reducing the need for manual intervention. By automating model selection, feature engineering, and hyperparameter tuning, AutoML can significantly enhance efficiency and adaptability in financial modeling, ensuring models remain relevant amidst rapidly changing market conditions.

*Explanability.* Despite their superior predictive capabilities, deep learning models often lack transparency, making them less preferable in domains where understanding decision-making processes is crucial, such as finance. The importance of risk management in finance cannot be overstated, necessitating models not just for their predictive performance but also for their ability to provide insights into their predictions. Developing task-specific, explainable AI methods is essential to demystify the "black box" of deep learning, offering clear insights into model decisions and fostering trust among stakeholders.

*Knowledge-driven AI.* Deep learning's reliance on extensive datasets poses challenges in investment scenarios characterized by sparse data, such as long-term value investing. In these contexts, knowledge-driven AI can offer a valuable complement to data-driven approaches, integrating domain expertise as a robust prior to compensate for the lack of large datasets.

*End-to-end modeling.* While the pipeline consists of multiple stages, in practice models for these stages are usually trained separately with varying goals. such inconsistency may hinder a coherent optimization direction. Hence it is also promising to explore a fully end-to-end modeling method that takes raw data and generate trades directly, while being trained directly from the final return. In this way the goals are aligned and may lead to better training results.

## 4 LARGE LANGUAGE MODELS FOR ALPHA RESEARCH

The rapid development of large language models has taken quantitative investment a significant step forward, from AI-powered to AI-automated, with vast prospects for future development. However, their application remains less mature compared to other deep learning methods, with limitations in certain sub-tasks. In this section, we categorize their roles into predictor and agent for alpha research, and explores other potential applications in subsection 4.3.

### 4.1 LLM-Based Predictor

The emergence of Large Language Models (LLMs) has redefined the role of textual data in quantitative finance, transforming how investors extract insights and generate predictive signals. Traditional models primarily relied on structured numerical data, but LLMs now enable a deeper understanding of financial narratives, capturing market sentiment, causal relationships, and latent trading

factors with unprecedented accuracy. This chapter explores how LLMs have evolved from passive sentiment classifiers to active market predictors, spanning applications from embedding-based sentiment extraction to direct factor generation for systematic trading. By bridging quantitative signals with qualitative reasoning, LLMs are paving the way for a new era of AI-driven market prediction—one where natural language understanding becomes a core component of financial modeling.
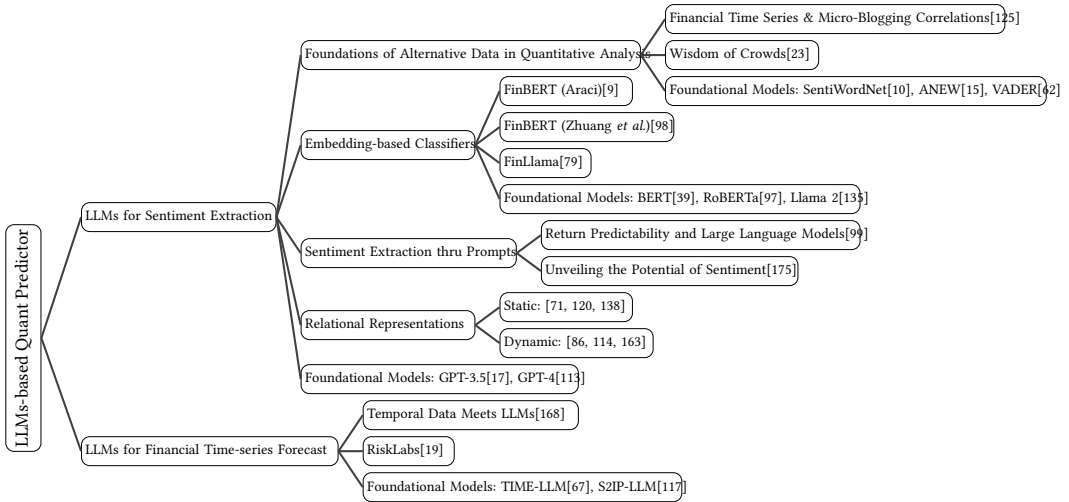


Fig. 7. LLM as predictors: an overview

### 4.1.1 LLMs for Sentiment Extraction.

*Foundations of Alternative Data Analysis in Financial Markets.* Traditional quantitative models predominantly rely on structured market data to generate trading signals. However, the rise of alternative data—particularly from digital financial activities and online discussions—has introduced new dimensions to market analysis. As investor sentiment and information dissemination accelerate through news articles, analytical reports, and social media, these textual data sources have become increasingly influential in shaping market movements.

The predictive potential of textual data in financial markets was recognized well before the advent of Large Language Models (LLMs) [23, 125]. Early approaches, however, lacked the computational sophistication needed to process large-scale textual information effectively. Initial methods relied on rudimentary techniques such as word frequency analysis and basic sentiment scoring [130]. For instance, [51] applied a Naïve Bayesian classifier to map financial news articles to stock movement labels based on predefined time intervals. While these early models demonstrated the feasibility of extracting signals from financial text, their effectiveness was constrained by simplistic feature representations and limited contextual understanding. Prior to LLMs, sentiment analysis tools such as SentiWordNet [10], ANEW [15], and VADER [62] were widely used, but they struggled with contextual ambiguity and lacked deep semantic comprehension.

As deep learning methodologies advanced, researchers sought to enhance sentiment analysis pipelines by integrating machine learning and neural networks. A notable example is the ensemble deep learning model proposed by [88], which combined sentiment scores from VADER with a hybrid recurrent neural network (RNN) approach. Their blending ensemble model improved stock

movement prediction accuracy by leveraging both textual news sentiment and stock price time-series data. A key contribution of their work was the exploration of historical news window sizes, reinforcing the notion that financial news exhibits long-memory effects, a characteristic later observed in LLM-driven sentiment models. These advancements paved the way for more context-aware and dynamic sentiment classification techniques, culminating in the emergence of LLM-powered financial text analysis.

Beyond sentiment scoring, recent work has expanded the role of causal relationships in financial forecasting. One significant advancement is the Causality-Guided Multi-Memory Interaction Network (CMIN)[100], which integrates financial text sentiment with causality-enhanced stock correlations to improve prediction accuracy. CMIN introduces a causal attention mechanism based on transfer entropy, ensuring that stock dependencies are not merely statistical artifacts but reflect directional information flows. The model's multi-memory interaction framework allows textual and market correlation features to reinforce each other dynamically, bridging a key gap between traditional sentiment models and structured financial indicators. By demonstrating that textual sentiment alone is insufficient without context-aware stock relationships, CMIN represents a broader evolution in alternative data analysis, moving toward more interpretable and multi-modal forecasting techniques.

The emergence of LLMs has fundamentally reshaped this landscape, enabling more nuanced and context-aware textual analysis. By leveraging deep learning architectures, these models can capture complex semantic relationships, sentiment nuances, and implicit market signals embedded in financial text. This advancement has significantly enhanced the predictive capabilities of text-based trading strategies, marking a paradigm shift in alternative data analysis for financial markets.

One of the most direct and impactful contributions of Large Language Models (LLMs) to quantitative finance is their ability to extract sentiment from financial texts. Traditional quantitative models primarily rely on structured numerical data, which inherently lacks the contextual and qualitative aspects embedded in financial news, analyst reports, earnings calls, and social media discussions. The integration of LLMs into sentiment extraction can be broadly categorized into two key approaches: *embedding-based classifiers* and *LLM-powered sentiment generation*.

*Embedding-Based LLM Classifiers.* One approach involves leveraging pre-trained word or sentence embeddings from models such as BERT[39], RoBERTa[97], or domain-specific variants like the two versions of FinBERT[9, 98]. These embeddings serve as dense vector representations of financial texts, capturing semantic relationships between words and phrases. By training supervised classifiers (e.g., logistic regression, random forests, or neural networks) on these embeddings, researchers can map textual sentiment to predefined categories—such as *bullish*, *bearish*, or *neutral*—with high precision. These classifiers benefit from efficiency and interpretability, making them ideal for large-scale applications in high-frequency trading and quantitative strategies. Further examples also include FinLlama[79] a variant of Llama-2-7B[135] fine-tuned on domain-specific financial text with a softmax activation layer for sentiment classification. In FinLlama, the fine-tuned model's sentiment classification is used directly as factor and back-tests were conducted to illustrate how more accurate sentiment classification can improve trading results. Another recent advancement in this direction is FININ[144], which extends traditional embedding-based sentiment classification by modeling interactions among financial news articles rather than treating each piece of news independently. FININ constructs a Financial Interconnected News Influence Network, integrating multi-modal data (news text embeddings and market statistics) to assess the impact of both individual news items and their contextual interactions. Notably, FININ leverages pre-trained LLMs but does not fine-tune them for sentiment classification. Instead, it processes these embeddings through an influence quantifier, incorporating financial theories to assess market impact. By

demonstrating that news interactions shape market sentiment, FININ represents an advancement in embedding-based models, bridging the gap between static embeddings and dynamic market-aware sentiment analysis.

*LLM-Powered Sentiment Classification.* An alternative, more dynamic approach involves directly prompting LLMs to generate sentiment classifications by incorporating financial news or earnings reports into structured queries. Unlike traditional methods that rely solely on pre-trained embeddings, this approach allows the LLM to contextually interpret sentiment within news text in relation to evolving market conditions. One of the most influential studies in this domain is [99], which systematically evaluated ChatGPT's ability to forecast stock price movements based on financial news sentiment. Their findings demonstrated that ChatGPT[? ]-generated sentiment signals exhibit predictive power, even after accounting for traditional asset pricing factors, with statistically significant results in backtesting. Additionally, they observed that ChatGPT's sentiment assessments closely aligned with analyst consensus and market reactions, suggesting that LLMs can approximate human-like financial reasoning. To validate this, their experimental design involved prompting ChatGPT to classify financial news headlines as bullish, bearish, or neutral, constructing sentiment-based trading signals, and evaluating their predictive strength through an asset pricing regression framework and trading simulations.

Building on this framework, [175] applied a similar methodology to Chinese financial news, testing its effectiveness in a trading simulation of Chinese A-shares, further demonstrating the adaptability of LLM-driven sentiment analysis across different markets. Meanwhile, [156] explored a hybrid approach, incorporating both structured quantitative data (e.g., stock prices and trading volume) and unstructured text from stock-related tweets to provide ChatGPT with a more comprehensive contextual understanding. Their work highlighted the potential for multimodal LLM-driven trading signals, combining sentiment cues with fundamental market data to enhance prediction accuracy.

*LLM-Based Relational Representation.* The integration of alternative data, such as news information, into relational frameworks (especially graph as Figure 5) has been shown to boost price prediction accuracy. Previously, graph relation construction was primarily based on real-world relationships, manual labeling, price correlations, or clustering algorithms. In contrast, LLMs' inherent comprehension and reasoning capabilities offer new perspectives for financial text understanding and market relationship modeling.

In recent studies, [71, 120, 138] proposed the construction of static financial relations and validated the feasibility of extracting information from financial texts. In practice, more dynamic models with temporal variations are more widely adopted. [114, 163] leverages information from video-audio and news sources respectively, establishing semantic relations. Based on this, predictions for market variables such as price, volatility and volume were made using a graph network approach. [86] propose FinDKG to extract global financial dynamic knowledge graph for risk tracking and investing with a better economic trend understanding.

### 4.1.2 LLMs for Direct Time-Series Forecasting.
The integration of Large Language Models (LLMs) into time-series forecasting has introduced a paradigm shift in predictive modeling, leveraging their ability to process multi-modal data, capture complex dependencies, and provide human-readable explanations. Traditional deep learning architectures have demonstrated varying degrees of success in forecasting time-series movements. However, these models often struggle with cross-sequence reasoning, multi-source data fusion, and explainability, limiting their interpretability and robustness in dynamic environments. Recent studies have explored how LLMs can address these challenges, enabling more context-aware, interpretable, and adaptive forecasting frameworks.

A foundational step in this direction is the demonstration that LLMs can perform time-series forecasting in a zero-shot manner, requiring no explicit training on numerical sequences. [53] shows that LLMs, when prompted effectively, can match or exceed purpose-built forecasting models, highlighting their ability to generalize patterns and model complex temporal dependencies without additional fine-tuning. Extending this idea, [67] introduces a reprogramming framework that aligns time-series data with the natural language capabilities of LLMs, enabling improved adaptability and performance across diverse forecasting tasks while keeping the backbone model intact. Further advancing this paradigm, [117] proposes S2IP-LLM, which aligns the semantic space of LLMs with time-series embeddings, enabling a contextualized prompt learning mechanism. These studies establish LLMs as generalizable time-series forecasters, setting the stage for their application in more specialized financial forecasting contexts.

Building on these foundational works, recent efforts have extended LLM applications to financial time-series forecasting, where market dynamics demand both predictive accuracy and interpretability. One of the pioneering studies in this domain, [168] explores how LLMs can be applied to financial time-series forecasting with an emphasis on explainability. The study focuses on NASDAQ-100 stock prediction and proposes a methodology that integrates zero-shot and few-shot inference using GPT-4 to predict price movements, instruction-based fine-tuning of Open LLaMA, demonstrating that fine-tuning a public LLM can yield competitive results, and multi-modal data integration, incorporating stock price time-series data, financial news, and company metadata. A key finding is that GPT-4's reasoning capabilities, particularly when prompted with Chain-of-Thought (CoT) prompting, significantly improve predictive accuracy over traditional models like ARMA-GARCH and gradient boosting trees. The study also underscores the ability of LLMs to generate human-readable justifications for their forecasts, making them more transparent and interpretable compared to black-box machine learning models.

Beyond standard asset price prediction, the RiskLabs framework[19] extends LLM applications to financial risk prediction, integrating multi-modal financial data to estimate market volatility and Value-at-Risk (VaR). This model processes earnings conference calls (ECCs), analyzing both textual transcripts and vocal features (e.g., tone, sentiment), market-related time-series data, capturing price movements over multiple time horizons, and contextual news data, aligning financial reports and media sentiment with earnings announcements. RiskLabs employs an LLM-powered multi-task learning approach, where different modules extract, encode, and fuse insights from structured (numerical time-series) and unstructured (text/audio) data sources. Experimental results demonstrate that the RiskLabs model outperforms traditional statistical and machine learning models in volatility forecasting, highlighting the potential of LLM-driven multi-modal forecasting.

The application of LLMs to financial time-series forecasting represents a significant leap forward in market prediction and risk assessment. By integrating multi-modal data sources, fine-tuning for financial reasoning, and leveraging structured inference techniques, LLMs enhance both accuracy and interpretability in predictive finance. Future research should continue to explore hybrid approaches, combining LLM-powered reasoning with domain-specific econometric models, to further refine adaptive and explainable financial forecasting methodologies.

## 4.2 LLM-Based Quant Agent

The rapid development of large language models (LLMs) has catalyzed a paradigm shift in financial AI systems. Domain-specific financial LLMs like BloombergGPT [151], FinGPT [164], and PIXIU [157] in English markets, along with CFGPT [83], [181], [152] and DISC-FinLLM [25] in Chinese contexts, have demonstrated superior reasoning capabilities through pre-training on massive financial corpora and instruction-tuning on specialized tasks. These models have enabled the

creation of standardized benchmarks for evaluating financial reasoning, sentiment analysis, and decision-making capabilities.

However, standalone LLMs face inherent limitations in real-time financial applications due to temporal latency and numerical processing constraints. This challenge has spurred the emergence of LLM-powered quant agents - AI systems that combine linguistic reasoning with tool invocation capabilities to process dynamic market data. These agents operate as autonomous entities capable of environmental perception, multi-step decision-making, and action execution through API integrations.
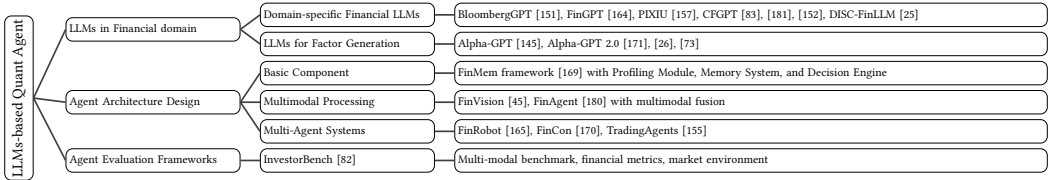


Fig. 8. Taxonomy of research in Financial LLMs and Agent Architectures.

*4.2.1 LLMs for Direct Factor Generation.* While the previous sections focused on LLM-powered sentiment extraction from financial alternative data, the application of Large Language Models in quantitative investment extends beyond it. LLMs can also be leveraged to generate predictive factors directly as a factor agent, offering a novel approach to feature engineering in trading models. Instead of merely analyzing sentiment, these models can extract latent patterns from unstructured financial data, summarize qualitative insights into numerical indicators, and integrate LLM-generated signals into systematic trading strategies. One of the most notable advancements in this area is Alpha-GPT [145], which introduces a human-AI interactive framework for factor mining in quantitative investment. Alpha-GPT leverages LLMs to assist researchers and traders in discovering novel alpha factors by engaging in an iterative dialogue, where the model proposes factor ideas, refines them based on human feedback, and generates executable code for implementation. This interactive workflow enables a more dynamic and adaptive approach to factor discovery, allowing domain experts to guide the AI's reasoning while benefiting from its vast knowledge and pattern recognition capabilities. By incorporating domain knowledge and real-time market context, Alpha-GPT enhances the traditional factor discovery process, reducing reliance on purely statistical factor mining techniques. In subsequent versions, Alpha-GPT 2.0 emphasizes the iterative and interactive factor analysis process between humans and AI [171]. Leveraging LLMs, it automates the entire pipeline from alpha mining to modeling and analysis. Building on this theme of LLM-driven factor generation, [26] explores how ChatGPT and GPT-4[113]can be positioned as a surrogate financial analyst to generate novel stock return factors. Instead of providing direct access to financial data, ChatGPT is only informed of the data structure and schema, ensuring that the generated factors are derived purely from its knowledge base rather than existing factor models. Through prompt engineering, the model is tasked with conceptualizing innovative stock return factors based on fundamental trading attributes. ChatGPT then outputs Python code to compute these factors, which researchers validate for originality and novelty before backtesting. This approach highlights LLMs' capacity for creative financial factor discovery without relying on predefined quantitative models. Some of the other works in this realm also include [73] have attempted at using LLMs to analyze financial statement and predict a company's future earnings, they further proceeded to construct porfolios adjustment strategies based on such predictions and showed that it yielded high returns.

*4.2.2  Architecture of LLM-based Quant Agents.* As shown in Figure.9, a LLM-based Financial Agent integrates the three stages of financial decision-making into a cohesive process. The left part of Figure.9 details the data types used: fundamental (company profiles, financial reports), price-volume data collected from the data platform such as Yahoo Finance, text data(such as news, Bloomberg, Reddit and social messages from X.com), and multimedia (teleconferences, images, videos). These data feed into a Predictor Agent, which uses large language models to forecast stock trends (up or down).The right part of Figure.9 shows the subsequent stages. The Portfolio Optimization Agent uses stock predictions to allocate funds optimally, considering return objectives, constraints, and risk control. Finally, the Order Execution Agent implements these allocations by executing market orders efficiently, minimizing losses from market impacts. A feedback loop at the bottom indicates back-testing for refining the optimization model and workflow. This system integrates data analysis, prediction, optimization, and execution for robust, data-driven investment strategies.

In the model prediction phase, the agent leverages large language models to analyze diverse data sources, enhancing prediction accuracy through LLMs and advanced deep learning techniques. Currently, the majority of scholarly investigations have focused on this particular aspect[165, 169, 180]. These predictions then feed into the portfolio optimization stage, where the agent uses optimization algorithms to allocate assets, balancing risk and return based on investor preferences. Finally, in the order execution phase, the agent employs smart order routing and execution strategies to minimize market impact and maximize trade efficiency. However, both of these components still necessitate further in-depth research and exploration[155]. Throughout these stages, the LLM-based Financial Agent ensures data-driven, efficient, and adaptive decision-making, offering a significant edge in dynamic financial markets.
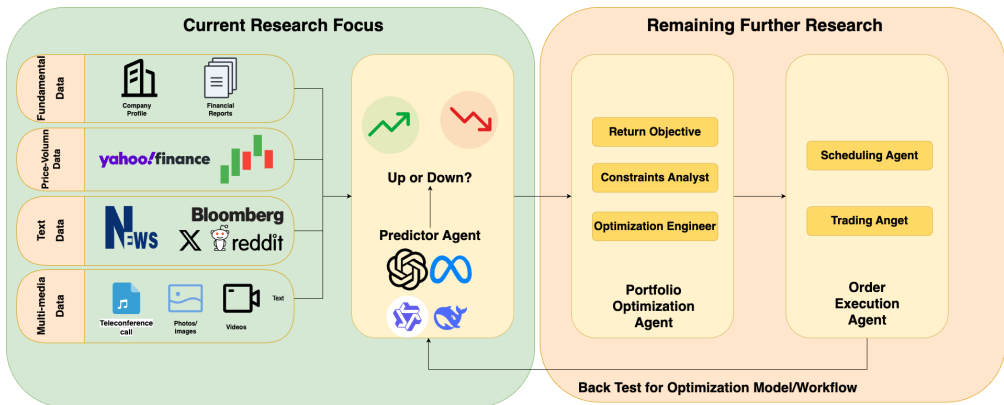


Fig. 9.  Architecture overview of LLM-based quant agents.It has three parts: using data to predict price trends, optimizing asset choices, and making trades.

*4.2.3  The Expansion of LLM-based Agents.* The evolution of financial agent architectures has demonstrated remarkable progress through successive innovations and expansions in cognitive modeling and system integration. The foundational FinMem framework [169] pioneered a basic architecture that established three core components now considered essential in modern systems. In the FINMEM framework, the design is structured around three core modules: The profiling module customizes the agent's characteristics by defining its professional background and investment risk inclinations, allowing it to adapt to market dynamics through dynamic risk settings. The memory module emulates human cognitive processes to handle hierarchical financial information,

incorporating working memory for temporary storage and operations, and layered long-term memory for managing information with varying time sensitivities. This structure enables the agent to prioritize and retrieve critical memory events effectively. The decision-making module integrates these outputs to support informed investment decisions, prioritizing key memory events and aligning with market conditions for high-quality trading outcomes. Experimental results demonstrate that FINMEM significantly outperforms other trading agents, achieving superior cumulative returns and Sharpe ratios, particularly in handling complex financial data and adapting to market conditions.

Subsequent expansions addressed the growing complexity of financial data ecosystems through multimodal processing architectures. The FinAgent [180] framework integrates multimodal data, including numerical, textual, and visual information, to analyze financial market dynamics and historical trading patterns. It employs a dual-level reflection module to adapt to market changes and enhance decision-making by learning from historical data. The agent utilizes large language models (LLMs) to process market intelligence and incorporates tool-augmented strategies to guide trading decisions. Experiments on six financial datasets, including stocks and cryptocurrencies, demonstrate that FinAgent significantly outperforms state-of-the-art baselines, showcasing its effectiveness in handling diverse financial assets and its ability to generalize across various market conditions. FinVision [45] is a multi-modal multi-agent framework designed for stock market prediction. The framework integrates specialized Large Language Model (LLM) agents to process various financial data types, including textual news summaries, technical analysis of candlestick charts, and reflection on historical trading signals. The experimental results demonstrate that FinVision outperforms traditional buy-and-hold strategies and reinforcement learning models in terms of annualized returns and risk-adjusted metrics, such as the Sharpe Ratio. While it falls short of the benchmark set by the FinAgent model, FinVision achieves competitive performance with a significantly shorter training period, highlighting its potential for adaptability and robustness in dynamic market conditions. The framework's ability to manage risk effectively while optimizing returns underscores its value in complex financial environments.

The architectural frontier has recently shifted towards collaborative multi-agent systems that replicate institutional trading desk dynamics. FinRobot [165] is an open-source AI agent platform designed for financial applications, leveraging large language models (LLMs) to enhance financial analysis and decision-making. The platform is structured into four layers: Financial AI Agents, Financial LLM Algorithms, LLMOps and DataOps, and Multi-source LLM Foundation Models. The Financial AI Agents layer uses Chain-of-Thought (CoT) prompting to break down complex financial problems into logical sequences, while the Financial LLM Algorithms layer dynamically configures model application strategies. The LLMOps and DataOps layer ensures accurate models through training and fine-tuning techniques, utilizing task-relevant data. The platform's Smart Scheduler mechanism integrates various LLMs, selecting the most suitable models for specific tasks. Experimental results demonstrate the platform's effectiveness in market forecasting and document analysis, showing its ability to provide comprehensive insights and actionable recommendations for financial professionals. FINCON [170] introduces a novel Large Language Model (LLM)-based multi-agent framework for financial decision-making, focusing on single-stock trading and portfolio management. The framework, inspired by real-world investment firm structures, employs a manager-analyst hierarchy to facilitate synchronized collaboration and reduce communication costs. It includes a dual-level risk-control component that monitors market risk and updates investment beliefs through self-critique. Experiments demonstrate that FINCON outperforms state-of-the-art LLM-based and DRL-based agents in terms of Cumulative Return (CR) and Sharpe Ratio (SR), while achieving one of the lowest Maximum Drawdown (MDD) values. Specifically, FINCON achieves a CR of 82.871% and a Sharpe Ratio of 1.972% for single-stock trading, and a CR of 113.836% and a

Sharpe Ratio of 3.269% for portfolio management, showcasing its robustness and effectiveness in managing financial risks and enhancing trading performance. The TradingAgents [155] framework introduces a novel multi-agent system for financial trading, leveraging Large Language Models (LLMs) to simulate the collaborative dynamics of a trading firm. The framework features specialized agents, including fundamental analysts, sentiment analysts, technical analysts, and traders, who collectively analyze market data and make informed trading decisions. Through structured communication protocols, agents engage in debates and discussions to reach balanced recommendations. Experimental results demonstrate that TradingAgents significantly outperforms baseline models in terms of cumulative returns, Sharpe ratio, and maximum drawdown, highlighting its ability to capture high returns while managing risk effectively. The framework's adaptability and explainability, facilitated by natural language-based operations, provide a distinct advantage over traditional trading strategies.

The maturation of financial agent research has driven systematic development of evaluation methodologies to address the domain's unique challenges. InvestorBench [82] introduces a comprehensive benchmark for evaluating large language model (LLM)-based agents in financial decision-making tasks. The benchmark addresses the lack of a versatile framework and standardized datasets by providing a suite of tasks applicable to various financial products, including stocks, cryptocurrencies, and exchange-traded funds (ETFs). The method involves a structured LLM agent framework with modules for perception, profile, memory, and action, which processes and integrates market data, historical insights, and self-reflection to inform investment decisions. Experimental results show that proprietary LLMs generally outperform open-source and domain-specific models, particularly in volatile market conditions. Larger model sizes within the open-source category also demonstrate improved performance, highlighting the importance of model scale in financial decision-making. The benchmark provides a valuable platform for assessing and comparing the reasoning and decision-making capabilities of LLMs in complex financial scenarios.

### 4.3 Current Limitations and Future Directions

*4.3.1 LLM-Based Predictor.* While the application of LLMs in financial sentiment analysis and time-series forecasting has shown promise, several critical challenges remain. Addressing these limitations will require novel adaptations and hybrid approaches that blend LLMs with specialized financial models to improve interpretability, efficiency, and predictive accuracy.

One key challenge in sentiment-based financial predictions is the misalignment between linguistic sentiment and market sentiment. LLMs, trained primarily on textual corpora, interpret sentiment through linguistic cues rather than financial market reactions. Although textual sentiment can correlate with price movements, it does not always translate directly into market behavior due to factors like pre-existing investor expectations, macroeconomic conditions, or sector-wide trends. Moving forward, future research should focus on aligning LLM-derived sentiment with financial market dynamics, potentially through reinforcement learning, cross-modal attention mechanisms, or fine-tuning on financial-specific datasets. By integrating structured financial indicators—such as historical price action, volatility metrics, and earnings reports—into the sentiment extraction pipeline, LLMs could develop a more market-aware understanding of sentiment beyond purely linguistic interpretations.

Another key limitation stems from the representation of numerical time-series data within LLMs. Because LLMs process text-based tokens rather than continuous numerical sequences, financial data must often be reformatted to align with an LLM's reasoning framework. However, time-series data has an inherently low signal-to-noise ratio, meaning that naive tokenization may introduce excessive noise, reducing predictive reliability. At this stage, the direct application of LLMs in time-series forecasting remains at the preliminary research, with limited adoption in the industry.

How to leverage the reasoning capabilities of LLMs for forecasting remains one of the key research focuses. A promising direction is to develop more structured embeddings for numerical time-series data, possibly by integrating transformer-based numerical encoders or creating specialized financial LLM tokenization techniques that preserve key statistical properties of market signals.

Similarly, in quant prediction, one of the most pressing concerns is latency in real-time decision-making. While zero-shot prompting and in-context learning enable LLMs to generalize across different forecasting tasks, these methods often introduce computational overhead. In fast-moving financial markets, models must generate predictions within milliseconds to remain actionable. Future advancements should explore lightweight, fine-tuned LLM architectures or hybrid models that combine LLM reasoning with low-latency numerical models (such as Kalman filters or traditional time-series regressions) to enhance real-time responsiveness.

Moreover, financial time-series forecasting is fundamentally relational, as market movements are influenced not only by a single company's historical prices but also by the behavior of its competitors, industry peers, and macroeconomic factors. Traditional econometric models explicitly encode these dependencies, but LLMs, when directly used for sequential forecasting, currently lack a clear mechanism to capture cross-company relationships. Future research should explore graph-based learning techniques, where LLMs incorporate relational embeddings that encode inter-stock dependencies. Alternatively, multi-modal LLM architectures that combine structured financial graphs with textual market news could enhance spatial dependency modeling and lead to more holistic market predictions.

In sum, while LLMs introduce powerful new capabilities for financial sentiment analysis and time-series forecasting, their current limitations highlight the need for domain-specific adaptations. Future research should focus on integrating market-aware sentiment modeling, optimizing real-time efficiency, developing structured numerical embeddings, and incorporating relational financial knowledge to further refine LLM-driven forecasting methodologies. By bridging the gap between text-driven reasoning and structured financial modeling, LLMs could evolve into more robust, interpretable, and adaptive tools for quantitative investment.

*4.3.2  LLM-Based Agent.* The integration of LLM-based agents into professional investment workflows has revealed significant gaps when evaluated against institutional-grade alpha generation pipelines (Section 2.2). These limitations span the entire investment decision chain, from predictive analytics to execution dynamics, presenting both challenges and opportunities for future research.

The prediction paradigm represents a fundamental challenge for current LLM agents. Same as LLM-based predictors, despite remarkable natural language processing capabilities, LLMs exhibit notable deficiencies in quantitative reasoning tasks [45, 180]. The numerical reasoning gap manifests in their limited ability to recognize precise price-volume patterns compared to specialized quant deep learning models, particularly in detecting subtle technical indicators and regime transitions. The first approach involves integrating existing deep learning models as tools into a specialized LLM quant-analyst agent, using the agent as the control center for generating and optimizing prediction models, thereby enhancing the accuracy of the agent's predictive and decision-making capabilities. Furthermore, Emerging hybrid architectures that combine LLMs with quantile regression networks offer promising solutions, potentially enabling probabilistic forecasting while maintaining interpretability through natural language explanations of model outputs.

Portfolio optimization frameworks in current LLM-based systems reveal significant deviations from institutional practice. Currently, trading agents [169, 180] predominantly focus on the trading individual assets, while research on multi-asset trading and portfolio optimization remains notably underdeveloped. The risk constraint formulations employed by existing frameworks often rely on oversimplified models, failing to incorporate sophisticated risk measures such as Conditional

Value at Risk (CVaR) . This limitation is exacerbated by the widespread neglect of transaction cost modeling, with most systems assuming frictionless markets. The diversification mechanics in current multi-agent systems lack explicit constraints, potentially leading to concentration risks that would be unacceptable in institutional portfolio optimization. The most straightforward solution is, since existing portfolio optimization techniques are relatively mature, to integrate them as tools into a specialized LLM optimization agent, enabling the agent to invoke these methods, analyze the results, execute decisions, evaluate and iterate autonomously, thereby enhancing the overall optimization performance of the agent. Neuro-symbolic approaches also present a promising direction for bridging this gap, potentially enabling the translation of LLM-generated market theses into mathematically rigorous optimization constraints while maintaining the flexibility of machine learning models.

The order execution stage exposes critical operational limitations in current agent designs. Most systems operate under the unrealistic assumption of perfect liquidity, ignoring the complex dynamics of limit order books and the substantial market impact of large orders [169, 170, 180]. This oversight is particularly problematic for liquid instruments where execution cost can significantly impact overall performance. The lack of real-time inference and adjustment mechanisms further compounds these issues, as agents fail to adapt their execution in time to changing market conditions. Integrating advanced order execution tools and market microstructure simulators with low-latency LLM reasoning modules could potentially enabling llm trading agents to develop sophisticated execution strategies that account for both market impact and opportunity cost.

The evolution of LLM-based quantitative agents must address these limitations through a holistic approach that considers the entire investment pipeline. Future research directions should focus on developing hybrid architectures that combine the strengths of LLMs in natural language processing and pattern recognition with the quantitative rigor of traditional financial models. This integration should span all stages of the investment process, from predictive analytics that incorporate both fundamental and technical factors, through portfolio construction frameworks that implement institution-grade risk management, to execution systems that account for real-world market frictions. Such advancements would bridge the gap between academic research and professional practice, potentially enabling the development of next-generation quantitative investment systems that combine the flexibility of machine learning with the robustness of traditional financial engineering.

## 5  CONCLUSION

Quantitative investment, particularly alpha strategy, as a forefront technology in financial markets, is attracting increasing attention. In this survey, we provide an in-depth and comprehensive review of how deep learning and large language models (LLMs) are being studied and applied in this domain. We present a thorough coverage of alpha strategy research, including data, models, and each components of the overall pipeline. Building on this foundation, we examine the transformative impact and performance improvements brought by deep learning in various aspects. In addition, we explore how LLMs can be utilized most effectively as predictors and agents. Finally, we compare different stages of development and their applications in quantitative investment models, summarizing the current limitations, analyzing the prevailing challenges and intricacies, and discussing a series of future research perspectives.

# REFERENCES

[1] 2022. Candlestick chart. https://en.wikipedia.org/w/index.php?title=Candlestick_chart&oldid=1091516507 Page Version ID: 1091516507.

[2] 2022. Factor investing. https://en.wikipedia.org/w/index.php?title=Factor_investing&oldid=1074680972 Page Version ID: 1074680972.

[3] 2022. Feature engineering. https://en.wikipedia.org/w/index.php?title=Feature_engineering&oldid=1090409817 Page Version ID: 1090409817.

[4] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[5] Almgren and Robert. 2012. Optimal Trading with Stochastic Liquidity and Volatility. *Siam Journal on Financial Mathematics* 3, 1 (2012), 163–181.

[6] Almgren and F. Robert. 2003. Optimal execution with nonlinear impact functions and trading-enhanced risk. *Applied Mathematical Finance* 10, 1 (2003), 1–18.

[7] Robert Almgren and Neil Chriss. 2000. Optimal execution of portfolio transactions. *Journal of Risk* (2000), 5–39.

[8] Hassan T Anis and Roy H Kwon. 2025. End-to-end, decision-based, cardinality-constrained portfolio optimization. *European Journal of Operational Research* 320, 3 (2025), 739–753.

[9] Dogu Araci. 2019. FinBERT: Financial Sentiment Analysis with Pre-trained Language Models. *arXiv preprint* (2019). arXiv:1908.10063

[10] Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. 2010. SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)* (Valletta, Malta). 2200–2204.

[11] Dirk Becherer, Todor Bilarev, and Peter Frentrup. 2016. Optimal Liquidation under Stochastic Liquidity. (2016).

[12] Dimitris Bertsimas and Andrew Lo. 1998. Optimal control of execution costs. *Journal of Financial Markets* 1, 1 (1998), 1–50. https://econpapers.repec.org/article/eeefinmar/v_3a1_3ay_3a1998_3ai_3a1_3ap_3a1-50.htm Publisher: Elsevier.

[13] Fischer Black. 1986. Noise. *The Journal of Finance* 41, 3 (1986), 528–543. https://doi.org/10.1111/j.1540-6261.1986.tb04513.x _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1540-6261.1986.tb04513.x.

[14] Anthony Brabazon, Michael Kampouridis, and Michael O'Neill. 2020. Applications of genetic programming to finance and economics: past, present, future. *Genetic Programming and Evolvable Machines* 21 (2020), 33–53.

[15] Margaret M. Bradley and Peter J. Lang. 1999. *Affective Norms for English Words (ANEW): Instruction Manual and Affective Ratings*. Technical Report. Citeseer.

[16] Antonio Briola, Jeremy Turiel, and Tomaso Aste. 2020. *Deep Learning modeling of Limit Order Book: a comparative perspective*. Technical Report arXiv:2007.07319. arXiv. https://doi.org/10.48550/arXiv.2007.07319 arXiv:2007.07319 [cs, q-fin, stat] type: article.

[17] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. *Advances in Neural Information Processing Systems* 33 (2020), 1877–1901.

[18] David Byrd, Maria Hybinette, and Tucker Hybinette Balch. 2020. ABIDES: Towards High-Fidelity Multi-Agent Market Simulation. In *Proceedings of the 2020 ACM SIGSIM Conference on Principles of Advanced Discrete Simulation*.

[19] Yupeng Cao, Zhi Chen, Qingyun Pei, Fabrizio Dimino, Lorenzo Ausiello, Prashant Kumar, K.P. Subbalakshmi, and Papa Momar Ndiaye. 2024. RiskLabs: Predicting Financial Risk Using Large Language Model Based on Multi-Sources Data. *arXiv preprint* arXiv:2404.07452 (2024). arXiv:2404.07452 [q-fin.RM] https://arxiv.org/abs/2404.07452

[20] Zhiyuan Cao, Jiayu Xu, Chengqi Dong, Peiwen Yu, and Tian Bai. 2024. MATCC: A Novel Approach for Robust Stock Price Prediction Incorporating Market Trends and Cross-time Correlations. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management* (Boise, ID, USA) *(CIKM '24)*. Association for Computing Machinery, New York, NY, USA, 187–196. https://doi.org/10.1145/3627673.3679715

[21] Marine Carrasco and Nérée Noumon. 2011. Optimal portfolio selection using regularization. *Citeseer, Tech. Rep.* (2011).

[22] Benjamin Remy Chabot, Eric Ghysels, and Ravi Jagannathan. 2014. Momentum Trading, Return Chasing, and Predictable Crashes. https://papers.ssrn.com/abstract=2521455

[23] Hailiang Chen, Prabuddha De, Yu Hu, and Byoung-Hyoun Hwang. 2013. Wisdom of Crowds: The Value of Stock Opinions Transmitted Through Social Media. *Review of Financial Studies* (2013).

[24] Sheng Chen and Hongxiang He. 2018. Stock prediction using convolutional neural network. In *IOP Conference series: materials science and engineering*, Vol. 435. IOP Publishing, 012026.

[25] Wei Chen, Qiushi Wang, Zefei Long, Xianyin Zhang, Zhongtian Lu, Bingxuan Li, Siyuan Wang, Jiarong Xu, Xiang Bai, Xuanjing Huang, et al. 2023. DISC-FinLLM: A Chinese financial large language model based on multiple experts fine-tuning. *arXiv preprint arXiv:2310.15205* (2023).

[26] Yifan Cheng and Kai Tang. 2024. GPT's Idea of Stock Factors. *Quantitative Finance* 24, 9 (2024), 1301–1326. https://doi.org/10.1080/14697688.2024.2318220

[27] Carl Chiarella and Giulia Iori. 2002. A simulation analysis of the microstructure of double auction markets. *Quantitative finance* 2, 5 (2002), 346.

[28] Carl Chiarella, Giulia Iori, and Josep Perello. 2009. The Impact of Heterogeneous Trading Rules on the Limit Order Book and Order Flows. *Journal of Economic Dynamics and Control* 33, 3 (2009), 525–537.

[29] Andrea Coletta, Joseph Jerome, Rahul Savani, and Svitlana Vyetrenko. 2023. Conditional Generators for Limit Order Book Environments: Explainability, Challenges, and Robustness. *CoRR* abs/2306.12806 (2023).

[30] Andrea Coletta, Aymeric Moulin, Svitlana Vyetrenko, and Tucker Balch. 2022. Learning to simulate realistic limit order book markets from data as a World Agent. In *Proceedings of the ACM International Conference on AI in Finance*.

[31] Robert T Collins, Alan J Lipton, Takeo Kanade, Hironobu Fujiyoshi, David Duggins, Yanghai Tsin, David Tolliver, Nobuyoshi Enomoto, Osamu Hasegawa, Peter Burt, et al. 2000. A system for video surveillance and monitoring. *VSAM final report* 2000, 1-68 (2000), 1.

[32] Rama Cont, Mihai Cucuringu, Renyuan Xu, and Chao Zhang. 2022. Tail-gan: Learning to simulate tail risk scenarios. *arXiv preprint arXiv:2203.01664* (2022).

[33] Stefania Corsaro and Valentina De Simone. 2019. Adaptive l 1-regularization for short-selling control in portfolio selection. *Computational Optimization and Applications* 72, 2 (2019), 457–478.

[34] Can Cui, Wei Wang, Meihui Zhang, Gang Chen, Zhaojing Luo, and Beng Chin Ooi. 2021. Alphaevolve: A learning framework to discover novel alphas in quantitative investment. In *Proceedings of the 2021 International conference on management of data*. 2208–2216.

[35] Kevin Dabérius, Elvin Granat, and Patrik Karlsson. 2019. Deep Execution - Value and Policy Based Reinforcement Learning for Trading and Beating Market Benchmarks. https://doi.org/10.2139/ssrn.3374766

[36] Divyanshu Daiya, Monika Yadav, and Harshit Singh Rao. 2024. Diffstock: Probabilistic relational stock market predictions using diffusion models. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 7335–7339.

[37] Sanmay Das. 2001. Filters, Wrappers and a Boosting-Based Hybrid for Feature Selection. In *Proceedings of the Eighteenth International Conference on Machine Learning (ICML '01)*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 74–81.

[38] Jacob Devlin. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).

[39] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*. https://api.semanticscholar.org/CorpusID:52967399

[40] Ryan Donnelly. 2022. Optimal execution: A review. *Applied Mathematical Finance* 29, 3 (2022), 181–212.

[41] Kelvin Du, Rui Mao, Frank Xing, and Erik Cambria. 2024. Explainable stock price movement prediction using contrastive learning. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*. 529–537.

[42] Yitong Duan, Lei Wang, Qizhong Zhang, and Jian Li. 2022. Factorvae: A probabilistic dynamic factor model based on variational autoencoder for predicting cross-sectional stock returns. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 36. 4468–4476.

[43] Eugene F. Fama and Kenneth R. French. 1992. The Cross-Section of Expected Stock Returns. *The Journal of Finance* 47, 2 (1992), 427–465. https://doi.org/10.1111/j.1540-6261.1992.tb04398.x _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1540-6261.1992.tb04398.x.

[44] Yuchen Fang, Kan Ren, Weiqing Liu, Dong Zhou, Weinan Zhang, Jiang Bian, Yong Yu, and Tie-Yan Liu. 2021. Universal Trading for Order Execution with Oracle Policy Distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 107–115. https://doi.org/10.1609/aaai.v35i1.16083

[45] Sorouralsadat Fatemi and Yuheng Hu. 2024. FinVision: A Multi-Agent Framework for Stock Market Prediction. In *Proceedings of the 5th ACM International Conference on AI in Finance*. 582–590.

[46] Fuli Feng, Xiangnan He, Xiang Wang, Cheng Luo, Yiqun Liu, and Tat-Seng Chua. 2019. Temporal Relational Ranking for Stock Prediction. *ACM Transactions on Information Systems* 37, 2 (March 2019), 1–30. https://doi.org/10.1145/3309547 arXiv: 1809.09441.

[47] Joao Fonseca and Fernando Bacao. 2023. Tabular and latent space synthetic data generation: a literature review. *Journal of Big Data* 10, 1 (2023), 115.

[48] P. A. Forsyth, J. S. Kennedy, S. T. Tse, and H. Windcliff. 2012. Optimal trade execution: A mean quadratic variation approach. *Journal of Economic Dynamics & Control* 36, 12 (2012), 1971–1991.

[49] Jean Pierre Fouque, Sebastian Jaimungal, and Yuri F. Saporito. 2021. Optimal Trading with Signals and Stochastic Price Impact. *Papers* (2021).

[50] Rao Fu, Jie Chen, Shutian Zeng, Yiping Zhuang, and Agus Sudjianto. 2019. Time series simulation by conditional generative adversarial net. *arXiv preprint arXiv:1904.11419* (2019).

[51] Győző Gidófalvi. 2001. *Using News Articles to Predict Stock Price Movements.* Technical Report. University of California, San Diego, Department of Computer Science and Engineering.

[52] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. *Advances in neural information processing systems* 27 (2014).

[53] Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew Gordon Wilson. 2023. Large Language Models Are Zero-Shot Time Series Forecasters. *arXiv preprint arXiv:2310.07820* (2023). https://arxiv.org/abs/2310.07820

[54] Mark A. Hall. 1999. *Correlation-based Feature Selection for Machine Learning.* PhD Thesis.

[55] Ben Hambly, Renyuan Xu, and Huining Yang. 2020. Policy Gradient Methods for the Noisy Linear Quadratic Regulator over a Finite Horizon. (2020).

[56] Dieter Hendricks and Diane Wilcox. 2014. A reinforcement learning extension to the Almgren-Chriss model for optimal trade execution. *Papers* (2014), 457–464.

[57] Ehsan Hoseinzade and Saman Haratizadeh. 2019. CNNpred: CNN-based stock market prediction using a diverse set of variables. *Expert Systems with Applications* 129 (2019), 273–285.

[58] Xiurui Hou, Kai Wang, Cheng Zhong, and Zhi Wei. 2021. ST-Trader: A Spatial-Temporal Deep Neural Network for Modeling Stock Market Movement. *IEEE/CAA Journal of Automatica Sinica* 8, 5 (May 2021), 1015–1024. https://doi.org/10.1109/JAS.2021.1003976

[59] Ziniu Hu, Weiqing Liu, Jiang Bian, Xuanzhe Liu, and Tie-Yan Liu. 2019. Listening to Chaotic Whispers: A Deep Learning Framework for News-oriented Stock Trend Prediction. *arXiv:1712.02136 [cs, q-fin]* (Feb. 2019). http://arxiv.org/abs/1712.02136 arXiv: 1712.02136.

[60] Hongbin Huang, Minghua Chen, and Xiao Qiao. 2024. Generative Learning for Financial Time Series with Irregular and Scale-Invariant Patterns. In *The Twelfth International Conference on Learning Representations.*

[61] Yu-Hao Huang, Chang Xu, Yang Liu, Weiqing Liu, Wu-Jun Li, and Jiang Bian. 2024. Controllable Financial Market Generation with Diffusion Guided Meta Agent. arXiv:2408.12991 [cs.CE] https://arxiv.org/abs/2408.12991

[62] C.J. Hutto and Eric Gilbert. 2015. VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text. *Proceedings of the 8th International Conference on Weblogs and Social Media, ICWSM 2014.*

[63] Kentaro Imajo, Kentaro Minami, Katsuya Ito, and Kei Nakagawa. 2021. Deep Portfolio Optimization via Distributional Prediction of Residual Factors. *Proceedings of the AAAI Conference on Artificial Intelligence* 35, 1 (May 2021), 213–222. https://doi.org/10.1609/aaai.v35i1.16095 Number: 1.

[64] Vinesh Jha. 2018. Implementing Alternative Data in an Investment Process. In *Big Data and Machine Learning in Quantitative Investment.* John Wiley & Sons, Ltd, 51–73. https://doi.org/10.1002/9781119522225.ch4 Section: 4 _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/9781119522225.ch4.

[65] Zhengyao Jiang, Dixing Xu, and Jinjun Liang. 2017. A Deep Reinforcement Learning Framework for the Financial Portfolio Management Problem. https://doi.org/10.48550/arXiv.1706.10059 arXiv:1706.10059 [cs, q-fin].

[66] Xianfeng Jiao, Zizhong Li, Chang Xu, Yang Liu, Weiqing Liu, and Jiang Bian. 2023. Microstructure-Empowered Stock Factor Extraction and Utilization. *arXiv preprint arXiv:2308.08135* (2023).

[67] Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y. Zhang, Xiaoming Shi, Pin-Yu Chen, Yuxuan Liang, Yuan-Fang Li, Shirui Pan, and Qingsong Wen. 2024. TIME-LLM: Time Series Forecasting by Reprogramming Large Language Models. In *International Conference on Learning Representations (ICLR).* arXiv:2310.01728 [cs.LG] https://arxiv.org/abs/2310.01728

[68] Zhigang Jin, Yang Yang, and Yuhong Liu. 2020. Stock closing price prediction based on sentiment analysis and LSTM. *Neural Computing and Applications* 32 (2020), 9713–9729.

[69] Zura Kakushadze. 2016. 101 formulaic alphas. *Wilmott* 2016, 84 (2016), 72–81.

[70] N Kannan. 2024. A review of Deep Generative Models for Synthetic Financial Data Generation. *Journal ID* 1233 (2024), 1259.

[71] Simerjot Kaur, Charese Smiley, Akshat Gupta, Joy Sain, Dongsheng Wang, Suchetha Siddagangappa, Toyin Aguda, and Sameena Shah. 2023. REFinD: Relation extraction financial dataset. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval.* 3054–3063.

[72] Idris Kharroubi and Huyen Pham. 2010. Optimal portfolio liquidation with execution cost and risk. *SIAM Journal on Financial Mathematics* (2010).

[73] Alex Kim, Michael Muhn, and Valeri Nikolaev. 2024. Financial Statement Analysis with Large Language Models. *arXiv preprint* (Jul 2024). arXiv:2407.17866

[74] Raehyun Kim, ChanHo So, Minbyul Jeong, SangHoon Lee, JinKyu Kim, and Jaewoo Kang. 2019. HATS: A Hierarchical Graph Attention Network for Stock Movement Prediction. *Research Papers in Economics,Research Papers in Economics* (Aug 2019).

[75] Raehyun Kim, Chan Ho So, Minbyul Jeong, Sanghoon Lee, Jinkyu Kim, and Jaewoo Kang. 2019. HATS: A Hierarchical Graph Attention Network for Stock Movement Prediction. *arXiv:1908.07999 [cs, q-fin]* (Nov. 2019). http://arxiv.org/abs/1908.07999 arXiv: 1908.07999.

[76] Diederik P Kingma. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).

[77] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. *arXiv:1609.02907 [cs, stat]* (Feb. 2017). http://arxiv.org/abs/1609.02907 arXiv: 1609.02907.

[78] Daphne Koller and Mehran Sahami. 1996. Toward optimal feature selection. In *Proceedings of the Thirteenth International Conference on International Conference on Machine Learning (ICML'96)*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 284–292.

[79] Thanos Konstantinidis, Giorgos Iacovides, Mingxue Xu, Tony G Constantinides, and Danilo Mandic. 2024. FinLlama: Financial Sentiment Classification for Algorithmic Trading Applications. *arXiv preprint arXiv:2403.12285* (2024).

[80] Zhao-Rong Lai and Haisheng Yang. 2022. A Survey on Gaps between Mean-Variance Approach and Exponential Growth Rate Approach for Portfolio Optimization. *Comput. Surveys* 55, 2 (2022), 25:1–25:36. https://doi.org/10.1145/3485274

[81] Bin Li and Steven C. H. Hoi. 2014. Online portfolio selection: A survey. *Comput. Surveys* 46, 3 (Jan. 2014), 35:1–35:36. https://doi.org/10.1145/2512962

[82] Haohang Li, Yupeng Cao, Yangyang Yu, Shashidhar Reddy Javaji, Zhiyang Deng, Yueru He, Yuechen Jiang, Zining Zhu, Koduvayur Subbalakshmi, Guojun Xiong, et al. 2024. INVESTORBENCH: A Benchmark for Financial Decision-Making Tasks with LLM-based Agent. *arXiv preprint arXiv:2412.18174* (2024).

[83] Jiangtong Li, Yuxuan Bian, Guoxuan Wang, Yang Lei, Dawei Cheng, Zhijun Ding, and Changjun Jiang. 2023. CFGPT: Chinese financial assistant with large language model. *arXiv preprint arXiv:2309.10654* (2023).

[84] Junjie Li, Yang Liu, Weiqing Liu, Shikai Fang, Lewen Wang, Chang Xu, and Jiang Bian. 2024. MarS: a Financial Market Simulation Engine Powered by Generative Foundation Model. *arXiv preprint arXiv:2409.07486* (2024).

[85] Junyi Li, Xintong Wang, Yaoyang Lin, Arunesh Sinha, and Michael P. Wellman. 2020. Generating Realistic Stock Market Order Streams. In *Proceedings of the AAAI Conference on Artificial Intelligence, AAAI*.

[86] Victor Xiaohui Li. 2023. Findkg: Dynamic knowledge graph with large language models for global finance. *Available at SSRN 4608445* (2023).

[87] Wei Li, Ruihan Bao, Keiko Harimoto, Deli Chen, Jingjing Xu, and Qi Su. 2020. Modeling the Stock Relation with Graph Network for Overnight Stock Movement Prediction. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence.* International Joint Conferences on Artificial Intelligence Organization, Yokohama, Japan, 4541–4547. https://doi.org/10.24963/ijcai.2020/626

[88] Yang Li and Yi Pan. 2022. A Novel Ensemble Deep Learning Model for Stock Prediction Based on Stock Prices and News. *International Journal of Data Science and Analytics* 13, 2 (2022), 139–149. https://doi.org/10.1007/s41060-021-00261-x

[89] Siyu Lin and Peter A. Beling. 2020. A Deep Reinforcement Learning Framework for Optimal Trade Execution. In *Machine Learning and Knowledge Discovery in Databases. Applied Data Science and Demo Track - European Conference, ECML PKDD 2020, Ghent, Belgium, September 14-18, 2020, Proceedings, Part V (Lecture Notes in Computer Science, Vol. 12461)*, Yuxiao Dong, Georgiana Ifrim, Dunja Mladenic, Craig Saunders, and Sofie Van Hoecke (Eds.). Springer, 223–240. https://doi.org/10.1007/978-3-030-67670-4_14

[90] Siyu Lin and Peter A. Beling. 2020. An End-to-End Optimal Trade Execution Framework based on Proximal Policy Optimization, Vol. 5. 4548–4554. https://doi.org/10.24963/ijcai.2020/627 ISSN: 1045-0823.

[91] Siyu Lin and Peter A. Beling. 2021. An Agent-Based Market Simulator for Back-Testing Deep Reinforcement Learning Based Trade Execution Strategies. In *Neural Information Processing - 28th International Conference, ICONIP 2021, Sanur, Bali, Indonesia, December 8-12, 2021, Proceedings, Part III (Lecture Notes in Computer Science, Vol. 13110)*, Teddy Mantoro, Minho Lee, Media Anugerah Ayu, Kok Wai Wong, and Achmad Nizar Hidayanto (Eds.). Springer, 644–653. https://doi.org/10.1007/978-3-030-92238-2_53

[92] Siyu Lin and Peter A. Beling. 2021. Investigating the Robustness and Generalizability of Deep Reinforcement Learning Based Optimal Trade Execution Systems. In *Intelligent Computing - Proceedings of the 2021 Computing Conference, Volume 2, SAI 2021, Virtual Event, 15-16 July, 2021 (Lecture Notes in Networks and Systems, Vol. 284)*, Kohei Arai (Ed.). Springer, 912–926. https://doi.org/10.1007/978-3-030-80126-7_64

[93] H. Liu and H. Motoda. 1998. *Feature Extraction, Construction and Selection: A Data Mining Perspective.* Springer US. https://books.google.bj/books?id=zi_0EdWW5fYC

[94] Qi Liu and Yue Zhang. 2018. Mining Evidences for Concept Stock Recommendation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 2103–2112. https://doi.org/10.18653/v1/N18-1191

[95] Qi Liu, Yue Zhang, and Jiangming Liu. 2018. Learning Domain Representation for Multi-Domain Sentiment Classification. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 541–550. https://doi.org/10.18653/v1/N18-1050

[96] Tom Liu, Stephen Roberts, and Stefan Zohren. 2023. Deep Inception Networks: A General End-to-End Framework for Multi-asset Quantitative Strategies. *arXiv preprint arXiv:2307.05522* (2023).

[97] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint* (2019). arXiv:1907.11692

[98] Zhuang Liu, Degen Huang, Kaiyu Huang, Zhuang Li, and Jun Zhao. 2020. FinBERT: A Pre-trained Financial Language Representation Model for Financial Text Mining. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, Christian Bessiere (Ed.). International Joint Conferences on Artificial Intelligence Organization, 4513–4519. https://doi.org/10.24963/ijcai.2020/622 Special Track on AI in FinTech.

[99] Andres Lopez Lira and Youyi Tang. 2023. Can ChatGPT Forecast Stock Price Movements? Return Predictability and Large Language Models. *arXiv preprint* (2023). https://doi.org/10.48550/arXiv.2304.07619 arXiv:2304.07619

[100] Di Luo, Weiheng Liao, Shuqi Li, Xin Cheng, and Rui Yan. 2023. Causality-Guided Multi-Memory Interaction Network for Multivariate Stock Price Movement Prediction. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL)*. Association for Computational Linguistics, Toronto, Canada, 12164–12176. https://doi.org/10.18653/v1/2023.acl-long.679

[101] Kolemann Lutz. 2018. How Satellite Imagery is Revolutionizing the Way we Invest. https://kolemannlutz.medium.com/how-satellite-imagery-is-revolutionizing-the-wa-da44bfb95d89

[102] Thomas Lux and Michele Marchesi. 1999. Scaling and criticality in a stochastic multi-agent model of a financial market. *Nature* 397, 6719 (1999), 498–500.

[103] Yilin Ma, Ruizhu Han, and Weizhong Wang. 2021. Portfolio optimization with return prediction using deep learning and machine learning. *Expert Syst. Appl.* 165 (2021), 113973. https://doi.org/10.1016/j.eswa.2020.113973

[104] M. Magdon-Ismail, A. Nicholson, and Y.S. Abu-Mostafa. 1998. Financial markets: very noisy information processing. *Proc. IEEE* 86, 11 (1998), 2184–2195. https://doi.org/10.1109/5.726786 Conference Name: Proceedings of the IEEE.

[105] Daiki Matsunaga, Toyotaro Suzumura, and Toshihiro Takahashi. 2019. Exploring Graph Neural Networks for Stock Market Predictions with Rolling Window Analysis. *RePEc: Research Papers in Economics - RePEc,RePEc: Research Papers in Economics - RePEc* (Jan 2019).

[106] Jose Menchero, Lei Ji, et al. 2019. Portfolio optimization with noisy covariance matrices. *The Journal of Investment Management* 17, 1 (2019), 200–206.

[107] Yuriy Nevmyvaka, Yi Feng, and Michael Kearns. 2006. Reinforcement learning for optimized trade execution. In *Machine Learning, Proceedings of the Twenty-Third International Conference (ICML 2006), Pittsburgh, Pennsylvania, USA, June 25-29, 2006*.

[108] Mahla Nikou, Gholamreza Mansourfar, and Jamshid Bagherzadeh. 2019. Stock price prediction using DEEP learning algorithm and its comparison with machine learning algorithms. *Intelligent Systems in Accounting, Finance and Management* 26, 4 (2019), 164–174.

[109] Brian Ning, Franco Ho Ting Lin, and Sebastian Jaimungal. 2021. Double Deep Q-Learning for Optimal Execution. *Applied Mathematical Finance* 28, 4 (July 2021), 361–380. https://doi.org/10.1080/1350486X.2022.2077783 Publisher: Routledge _eprint: https://doi.org/10.1080/1350486X.2022.2077783.

[110] Hui Niu, Siyuan Li, and Jian Li. 2022. MetaTrader: An Reinforcement Learning Approach Integrating Diverse Policies for Portfolio Optimization. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022*, Mohammad Al Hasan and Li Xiong (Eds.). ACM, 1573–1583. https://doi.org/10.1145/3511808.3557363

[111] Pisut Oncharoen and Peerapon Vateekul. 2018. Deep learning using risk-reward function for stock market prediction. In *Proceedings of the 2018 2nd International Conference on Computer Science and Artificial Intelligence*. 556–561.

[112] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. 2016. WaveNet: A Generative Model for Raw Audio. In *The 9th ISCA Speech Synthesis Workshop, Sunnyvale, CA, USA, 13-15 September 2016*. ISCA, 125. http://www.isca-speech.org/archive/SSW_2016/abstracts/ssw9_DS-4_van_den_Oord.html

[113] OpenAI. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774* (2023). https://arxiv.org/abs/2303.08774

[114] Kun Ouyang, Yi Liu, Shicheng Li, Ruihan Bao, Keiko Harimoto, and Xu Sun. 2024. Modal-adaptive knowledge-enhanced graph-based financial prediction from monetary policy conference calls with LLM. *arXiv preprint arXiv:2403.16055* (2024).

[115] Szilard Pafka and Imre Kondor. 2004. Estimated correlation matrices and portfolio optimization. *Physica A: statistical mechanics and its applications* 343 (2004), 623–634.

[116] R.G. Palmer, W. Brian Arthur, John H. Holland, Blake LeBaron, and Paul Tayler. 1994. Artificial economic life: a simple model of a stockmarket. *Physica D: Nonlinear Phenomena* 75, 1 (1994), 264–274.

[117] Zijie Pan, Yushan Jiang, Sahil Garg, Anderson Schneider, Yuriy Nevmyvaka, and Dongjin Song. 2024. S2IP-LLM: Semantic Space Informed Prompt Learning with LLM for Time Series Forecasting. *arXiv preprint* arXiv:2403.05798 (2024). arXiv:2403.05798 [cs.LG] https://arxiv.org/abs/2403.05798

[118] Frank Partnoy. 2019. Stock Picks From Space. https://www.theatlantic.com/magazine/archive/2019/05/stock-value-satellite-images-investing/586009/ Section: Business.

[119] Marco Raberto, Silvano Cincotti, Sergio M. Focardi, and Michele Marchesi. 2001. Agent-based simulation of a financial market. *Physica A: Statistical Mechanics and its Applications* 299, 1 (2001), 319–327.

[120] Pawan Kumar Rajpoot and Ankur Parikh. 2023. Gpt-finre: in-context learning for financial relation extraction using large language models. *arXiv preprint arXiv:2306.17519* (2023).

[121] Weizhe Ren, Yichen Qin, and Yang Li. 2024. Alpha Mining and Enhancing via Warm Start Genetic Programming for Quantitative Investment. *arXiv preprint arXiv:2412.00896* (2024).

[122] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. 2014. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*. PMLR, 1278–1286.

[123] Alan Ritter, Oren Etzioni, and Sam Clark. 2012. Open domain event extraction from twitter. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. 1104–1112.

[124] Stephen A Ross. 1976. The arbitrage theory of capital asset pricing. *Journal of Economic Theory* 13, 3 (Dec. 1976), 341–360. https://doi.org/10.1016/0022-0531(76)90046-6

[125] Eduardo J. Ruiz, Vagelis Hristidis, Carlos Castillo, Aristides Gionis, and Alejandro Jaimes. 2012. Correlating Financial Time Series with Micro-blogging Activity. In *Proceedings of the Fifth ACM International Conference on Web Search and Data Mining (WSDM)* (Seattle, Washington, USA). ACM, New York, NY, USA, 513–522.

[126] Timur Sattarov, Marco Schreyer, and Damian Borth. 2023. Findiff: Diffusion models for financial tabular data generation. In *Proceedings of the Fourth ACM International Conference on AI in Finance*. 64–72.

[127] Ramit Sawhney, Shivam Agarwal, Arnav Wadhwa, and Rajiv Ratn Shah. 2020. Spatiotemporal Hypergraph Convolution Network for Stock Movement Forecasting. In *2020 IEEE International Conference on Data Mining (ICDM)*. 482–491. https://doi.org/10.1109/ICDM50108.2020.00057 ISSN: 2374-8486.

[128] Ramit Sawhney, Arnav Wadhwa, Shivam Agarwal, and Rajiv Ratn Shah. 2021. FAST: Financial News and Tweet Based Time Aware Network for Stock Trading. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Association for Computational Linguistics, Online, 2164–2175. https://doi.org/10.18653/v1/2021.eacl-main.185

[129] Anatoly B Schmidt. 2019. Managing portfolio diversity within the mean variance theory. *Annals of Operations Research* 282, 1 (2019), 315–329.

[130] Robert Schumaker and Hsinchun Chen. 2009. Textual Analysis of Stock Market Prediction Using Breaking Financial News: The AZFin Text System. *ACM Transactions on Information Systems* 27 (2009). https://doi.org/10.1145/1462198.1462204

[131] Omer Berat Sezer, Mehmet Ugur Gudelek, and Ahmet Murat Ozbayoglu. 2020. Financial time series forecasting with deep learning : A systematic literature review: 2005–2019. *Applied Soft Computing* 90 (May 2020), 106181. https://doi.org/10.1016/j.asoc.2020.106181

[132] Hyun Sik Sim, Hae In Kim, and Jae Joon Ahn. 2019. Is deep learning for image recognition applicable to stock market prediction? *Complexity* 2019, 1 (2019), 4324878.

[133] HMHS Surendra, HS Mohan, et al. 2017. A review of synthetic data generation methods for privacy preserving data publishing. *International Journal of Scientific & Technology Research* 6, 3 (2017), 95–101.

[134] Gattaiah Tadoori and Yakanna Guguloth. 2020. *An Introduction to Quantamental Investing*. SSRN Scholarly Paper 3704670. Social Science Research Network, Rochester, NY. https://doi.org/10.2139/ssrn.3704670

[135] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Sharan Batra, Pratik Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. *arXiv preprint* (2023). arXiv:2307.09288

[136] Igor Tulchinsky. 2019. *Finding Alphas: A quantitative approach to building trading strategies*. John Wiley & Sons.

[137] A Sinem Uysal, Xiaoyue Li, and John M Mulvey. 2024. End-to-end risk budgeting portfolio optimization with neural networks. *Annals of Operations Research* 339, 1 (2024), 397–426.

[138] Dimitrios Vamvourellis, Máté Tóth, Snigdha Bhagat, Dhruv Desai, Dhagash Mehta, and Stefano Pasquali. 2024. Company similarity using large language models. In *2024 IEEE Symposium on Computational Intelligence for Financial Engineering and Economics (CIFEr)*. IEEE, 1–9.

[139] Milena Vuletić, Felix Prenzel, and Mihai Cucuringu. 2024. Fin-gan: Forecasting and classifying financial time series via generative adversarial networks. *Quantitative Finance* 24, 2 (2024), 175–199.

[140] Heyuan Wang, Shun Li, Tengjiao Wang, and Jiayi Zheng. 2021. Hierarchical Adaptive Temporal-Relational Modeling for Stock Trend Prediction.. In *IJCAI*. 3691–3698.

[141] Haoran Wang and Xun Yu Zhou. 2019. Continuous-Time Mean-Variance Portfolio Optimization via Reinforcement Learning. *CoRR* abs/1904.11392 (2019). http://arxiv.org/abs/1904.11392 arXiv: 1904.11392.

[142] Jianian Wang, Sheng Zhang, Yanghua Xiao, and Rui Song. 2021. A review on graph neural network methods in financial applications. *arXiv preprint arXiv:2111.15367* (2021).

[143] Jingyuan Wang, Yang Zhang, Ke Tang, Junjie Wu, and Zhang Xiong. 2019. AlphaStock: A Buying-Winners-and-Selling-Losers Investment Strategy using Interpretable Deep Reinforcement Attention Networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, Anchorage AK USA, 1900–1908. https://doi.org/10.1145/3292500.3330647

[144] Mengyu Wang, Shay B. Cohen, and Tiejun Ma. 2024. Modeling News Interactions and Influence for Financial Market Prediction. *arXiv preprint* cs.CE (2024). arXiv:2410.10614 [cs.CE] https://arxiv.org/abs/2410.10614

[145] Shuang Wang, Haotian Yuan, Li Zhou, Lionel M. Ni, Heung-Yeung Shum, and Jun Guo. 2023. Alpha-GPT: Human-AI Interactive Alpha Mining for Quantitative Investment. *arXiv preprint* (2023). arXiv:2308.00016

[146] Xintong Wang and Michael P. Wellman. 2017. Spoofing the Limit Order Book: An Agent-Based Model. In *Proceedings of the Conference on Autonomous Agents and MultiAgent Systems*.

[147] Eric W. Weisstein. [n. d.]. Hyperedge. https://mathworld.wolfram.com/ Publisher: Wolfram Research, Inc..

[148] Qingsong Wen, Tian Zhou, Chaoli Zhang, Weiqi Chen, Ziqing Ma, Junchi Yan, and Liang Sun. 2022. *Transformers in Time Series: A Survey*. Technical Report arXiv:2202.07125. arXiv. https://doi.org/10.48550/arXiv.2202.07125 arXiv:2202.07125 [cs, eess, stat] type: article.

[149] Magnus Wiese, Robert Knobloch, Ralf Korn, and Peter Kretschmer. 2020. Quant GANs: Deep Generation of Financial Time Series. *Quantitative Finance* 20, 9 (Sept. 2020), 1419–1440. https://doi.org/10.1080/14697688.2020.1730426 arXiv: 1907.06673.

[150] Magnus Wiese, Robert Knobloch, Ralf Korn, and Peter Kretschmer. 2020. Quant GANs: deep generation of financial time series. *Quantitative Finance* 20, 9 (2020), 1419–1440.

[151] Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564* (2023).

[152] Xiaojun Wu, Junxi Liu, Huanyi Su, Zhouchi Lin, Yiyan Qi, Chengjin Xu, Jiajun Su, Jiajie Zhong, Fuwei Wang, Saizhuo Wang, et al. 2024. Golden Touchstone: A Comprehensive Bilingual Benchmark for Evaluating Financial Large Language Models. *arXiv preprint arXiv:2411.06272* (2024).

[153] Yufei Wu, Daniele Magazzeni, and Manuela Veloso. 2021. How Robust are Limit Order Book Representations under Data Perturbation?. In *ICML Workshop on Representation Learning for Finance and E-Commerce Applications*.

[154] Haochong Xia, Shuo Sun, Xinrun Wang, and Bo An. 2024. Market-GAN: Adding Control to Financial Market Data Generation with Semantic Context. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 15996–16004.

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

[156] Qianqian Xie et al. 2023. The Wall Street Neophyte: A Zero-Shot Analysis of ChatGPT Over Multimodal Stock Movement Prediction Challenges. *arXiv preprint* (2023). arXiv:2304.05351

[157] Qianqian Xie, Weiguang Han, Xiao Zhang, Yanzhao Lai, Min Peng, Alejandro Lopez-Lira, and Jimin Huang. 2023. Pixiu: A large language model, instruction data and evaluation benchmark for finance. *arXiv preprint arXiv:2306.05443* (2023).

[158] Wentao Xu, Weiqing Liu, Lewen Wang, Yingce Xia, Jiang Bian, Jian Yin, and Tie-Yan Liu. 2021. Hist: A graph-based framework for stock trend forecasting via mining concept-oriented shared information. *arXiv preprint arXiv:2110.13716* (2021).

[159] Wentao Xu, Weiqing Liu, Lewen Wang, Yingce Xia, Jiang Bian, Jian Yin, and Tie-Yan Liu. 2022. HIST: A Graph-based Framework for Stock Trend Forecasting via Mining Concept-Oriented Shared Information. *arXiv:2110.13716 [cs, q-fin]* (Jan. 2022). http://arxiv.org/abs/2110.13716 arXiv: 2110.13716.

[160] Wentao Xu, Weiqing Liu, Chang Xu, Jiang Bian, Jian Yin, and Tie-Yan Liu. 2021. REST: Relational Event-driven Stock Trend Forecasting. *Proceedings of the Web Conference 2021* (April 2021), 1–10. https://doi.org/10.1145/3442381.3450032 arXiv: 2102.07372.

[161] Wentao Xu, Weiqing Liu, Chang Xu, Jiang Bian, Jian Yin, and Tie-Yan Liu. 2021. Rest: Relational event-driven stock trend forecasting. In *Proceedings of the web conference 2021*. 1–10.

[162] Yumo Xu and Shay B. Cohen. 2018. Stock Movement Prediction from Tweets and Historical Prices. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Melbourne, Australia, 1970–1979. https://doi.org/10.18653/v1/P18-1183

[163] Zhiyu Xu, Yi Liu, Yuchi Wang, Ruihan Bao, Keiko Harimoto, and Xu Sun. 2025. Modeling Interactions Between Stocks Using LLM-Enhanced Graphs for Volume Prediction. In *Proceedings of the Joint Workshop of the 9th Financial Technology and Natural Language Processing (FinNLP), the 6th Financial Narrative Processing (FNP), and the 1st Workshop on Large Language Models for Finance and Legal (LLMFinLegal)*. 153–163.

[164] Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. 2023. Fingpt: Open-source financial large language models. *arXiv preprint arXiv:2306.06031* (2023).

[165] Hongyang Yang, Boyu Zhang, Neng Wang, Cheng Guo, Xiaoli Zhang, Likun Lin, Junlin Wang, Tianyu Zhou, Mao Guan, Runjia Zhang, et al. 2024. FinRobot: an open-source AI agent platform for financial applications using large language models. *arXiv preprint arXiv:2405.14767* (2024).

[166] Zinuo You, Pengju Zhang, Jin Zheng, and John Cartlidge. 2024. Multi-Relational Graph Diffusion Neural Network with Parallel Retention for Stock Trends Classification. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 6545–6549. https://doi.org/10.1109/ICASSP48485.2024.10447394

[167] Shuo Yu, Hongyan Xue, Xiang Ao, Feiyang Pan, Jia He, Dandan Tu, and Qing He. 2023. Generating synergistic formulaic alpha collections via reinforcement learning. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 5476–5486.

[168] Xinli Yu, Zheng Chen, Yuan Ling, Shujing Dong, Zongyi Liu, and Yanbin Lu. 2023. Temporal Data Meets LLM–Explainable Financial Time Series Forecasting. *arXiv preprint* arXiv:2306.11025 (2023). arXiv:2306.11025 [cs.LG] https://arxiv.org/abs/2306.11025

[169] Yangyang Yu, Haohang Li, Zhi Chen, Yuechen Jiang, Yang Li, Denghui Zhang, Rong Liu, Jordan W Suchow, and Khaldoun Khashanah. 2024. FinMem: A performance-enhanced LLM trading agent with layered memory and character design. In *Proceedings of the AAAI Symposium Series*, Vol. 3. 595–597.

[170] Yangyang Yu, Zhiyuan Yao, Haohang Li, Zhiyang Deng, Yuechen Jiang, Yupeng Cao, Zhi Chen, Jordan Suchow, Zhenyu Cui, Rong Liu, et al. 2025. Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making. *Advances in Neural Information Processing Systems* 37 (2025), 137010–137045.

[171] Hang Yuan, Saizhuo Wang, and Jian Guo. 2024. Alpha-GPT 2.0: Human-in-the-Loop AI for quantitative investment. *arXiv preprint arXiv:2402.09746* (2024).

[172] Liang Zeng, Lei Wang, Hui Niu, Ruchen Zhang, Ling Wang, and Jian Li. 2021. Trade when opportunity comes: price movement forecasting via locality-aware attention and iterative refinement labeling. *arXiv preprint arXiv:2107.11972* (2021).

[173] Chong Zhang, Xinyi Liu, Mingyu Jin, Zhongmou Zhang, Lingyao Li, Zhengting Wang, Wenyue Hua, Dong Shu, Suiyuan Zhu, Xiaobo Jin, et al. 2024. When ai meets finance (stockagent): Large language model-based stock trading in simulated real-world environments. *arXiv preprint arXiv:2407.18957* (2024).

[174] Feng Zhang, Ruite Guo, and Honggao Cao. 2020. Information Coefficient as a Performance Measure of Stock Selection Models. https://doi.org/10.48550/arXiv.2010.08601 arXiv:2010.08601 [q-fin, stat].

[175] Hao Zhang, Feng Hua, Chao Xu, Jun Guo, Haoxiang Kong, and Rui Zuo. 2023. Unveiling the Potential of Sentiment: Can Large Language Models Predict Chinese Stock Price Movements? *arXiv preprint* (2023). https://doi.org/10.48550/arXiv.2306.14222 arXiv:2306.14222

[176] Liheng Zhang, Charu Aggarwal, and Guo-Jun Qi. 2017. Stock price prediction via discovering multi-frequency trading patterns. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*. 2141–2149.

[177] Liheng Zhang, Charu Aggarwal, and Guo-Jun Qi. 2017. Stock Price Prediction via Discovering Multi-Frequency Trading Patterns. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, Halifax NS Canada, 2141–2149. https://doi.org/10.1145/3097983.3098117

[178] Qiuyue Zhang, Chao Qin, Yunfeng Zhang, Fangxun Bao, Caiming Zhang, and Peide Liu. 2022. Transformer-based attention network for stock movement prediction. *Expert Systems with Applications* 202 (2022), 117239.

[179] Tianping Zhang, Yuanqi Li, Yifei Jin, and Jian Li. 2020. Autoalpha: an efficient hierarchical evolutionary algorithm for mining alpha factors in quantitative investment. *arXiv preprint arXiv:2002.08245* (2020).

[180] Wentao Zhang, Lingxuan Zhao, Haochong Xia, Shuo Sun, Jiaze Sun, Molei Qin, Xinyi Li, Yuqing Zhao, Yilei Zhao, Xinyu Cai, et al. 2024. A multimodal foundation agent for financial trading: Tool-augmented, diversified, and generalist. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4314–4325.

[181] Xuanyu Zhang and Qing Yang. 2023. Xuanyuan 2.0: A large chinese financial chat model with hundreds of billions parameters. In *Proceedings of the 32nd ACM international conference on information and knowledge management*.

4435–4439.

[182] Yifan Zhang, Peilin Zhao, Qingyao Wu, Bin Li, Junzhou Huang, and Mingkui Tan. 2022. Cost-Sensitive Portfolio Selection via Deep Reinforcement Learning. *IEEE Trans. Knowl. Data Eng.* 34, 1 (2022), 236–248. https://doi.org/10.1109/TKDE.2020.2979700

[183] Yunfei Zhang, Zhihua Zhou, Junwei Liu, and Jianjuan Yuan. 2022. Data augmentation for improving heating load prediction of heating substation based on TimeGAN. *Energy* 260 (2022), 124919.

[184] Zihao Zhang, Stefan Zohren, and Stephen Roberts. 2019. Deeplob: Deep convolutional neural networks for limit order books. *IEEE Transactions on Signal Processing* 67, 11 (2019), 3001–3012.

[185] Kaiping Zheng, Wei Wang, Jinyang Gao, Kee Yuan Ngiam, Beng Chin Ooi, and Wei Luen James Yip. 2017. Capturing feature-level irregularity in disease progression modeling. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management.* 1579–1588.

[186] Yanqiao Zhu, Weizhi Xu, Jinghao Zhang, Qiang Liu, Shu Wu, and Liang Wang. 2021. Deep Graph Structure Learning for Robust Representations: A Survey. *arXiv:2103.03036 [cs]* (March 2021). http://arxiv.org/abs/2103.03036 arXiv:2103.03036.