**UEH University**

**UEH Institute of Innovation**

**UEH College of Technology and Design**

**Final Report AI Project:**

**Investment Asset Management Advisory Services**

**Student team:** Group 01

**Student name & ID:** Đinh Tấn Lộc - 31221020226

Trần Thị Kim Chi - 31221026465

Nguyễn Thị Thanh Hương - 31221025560

Lê Thị Bảo Ngọc - 31221021291

Nguyễn Văn Phát - 31221025527

**Academic supervisor:** Ph.D Nguyễn Thiên Bảo

*Ho Chi Minh City, November 5th, 2024.*

# MEMBERS OF GROUP

| Ordinal number | Student ID | Full name | Contribution (%) |
|---|---|---|---|
| 01 | 31221020226 | Dinh Tan Loc | 100% |
| 02 | 31221025560 | Nguyen Thi Thanh Huong | 100% |
| 03 | 31221021291 | Le Thi Bao Ngoc | 100% |
| 04 | 31221026465 | Tran Thi Kim Chi | 100% |
| 05 | 31221025527 | Nguyen Van Phat | 100% |

# ABSTRACT

This study aims to develop an AI-powered investment support application for the Vietnamese financial market, integrating a chatbot to enhance user experience. The application leverages advanced technologies such as Natural Language Processing (NLP), Machine Learning, and Deep Learning to analyze market data, predict stock prices, optimize portfolios, and provide personalized investment recommendations. The chatbot is designed to answer user queries, deliver real-time market information, and offer technical support. Preliminary results show that the application effectively improves investment performance and user satisfaction. In the context of Industry 4.0, AI technology is being widely applied, especially in finance, where precision and speed in data analysis are crucial. However, Vietnam's financial market has yet to fully optimize smart technologies, presenting a great opportunity to develop intelligent investment solutions that assist investors in making more informed decisions. This research is motivated by the potential of AI to enhance decision-making processes and increase participation in the financial market, contributing to sustainable and inclusive economic growth. The necessity of this application is evident when considering the lack of sophisticated investment tools in Vietnam, particularly those capable of complex data analysis, trend prediction, and timely recommendations. Many individual investors struggle with gathering and analyzing information to make decisions. By automating this process, the application can reduce errors, maximize returns, and promote transparency and fairness in the financial market. The application is built on machine learning models, using supervised and unsupervised algorithms to identify trading patterns, predict market fluctuations, and optimize portfolios. Deep learning algorithms process real-time market data to enhance stock price predictions, while NLP enables the chatbot to understand natural language, interact with users, answer questions, and provide real-time investment insights. The practical significance of this application lies in its ability to support individual investors in making better decisions while enhancing their financial literacy by providing accessible and useful information. Additionally, it contributes to financial inclusion by expanding investment opportunities to a broader range of users, including those with limited financial knowledge, thereby fostering the sustainable development of Vietnam's financial market. This study not only advances the development of intelligent investment solutions but also creates opportunities for investors to access cutting-edge technologies, enhancing the global competitiveness of Vietnam's financial market.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| Ordinal number | Symbol abbreviation | Explain the meaning of the word |
|:---:|:---:|:---|
| 01 | AI | Application Programming Interface |
| 02 | API | Artificial Intelligence |
| 03 | AR | Autoregressive |
| 04 | ARIMA | Autoregressive integrated moving average |
| 05 | MAE | Mean Absolute Error |
| 06 | ARIMA | Auto-Regressive Integrated Moving Average |
| 07 | DL | Deep learning |
| 08 | EUT | Expected Utility Theory |
| 09 | LLM | Large Language Model |
| 10 | LSTM | Long short term memory |
| 11 | LSTM | Long Short-Term Memory |
| 12 | MAS | Machine learning |
| 13 | ML | Mean Squared Error |

| Ordinal number | Symbol abbreviation | Explain the meaning of the word |
| --- | --- | --- |
| 14 | MSE | Multi-Agent System |
| 15 | NLP | Natural Language Processing |
| 16 | NLP | Natural Language Processing |
| 17 | OCR | Optical Character Recognition |
| 18 | $R^2$ | Retrieval-Augmented Generation |
| 19 | RAG | R-squared (coefficient of determination) |
| 20 | SQL | Structured Query Language |
| 21 | XGBoost | Extreme Gradient Boosting |

# THE REPORT

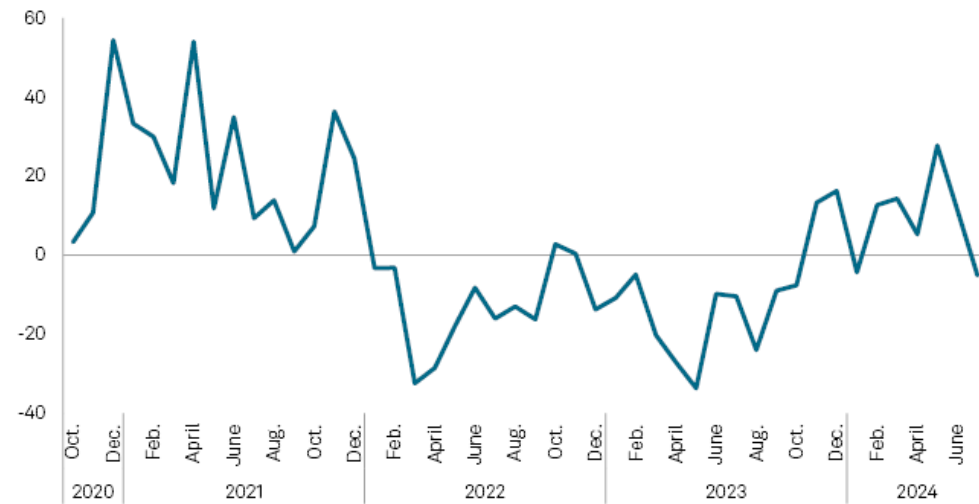## I. Introduction

## 1.1. Problem define

The stock market is an essential element of contemporary economies, providing both individuals and institutions with avenues for investment in various companies, which can lead to wealth accumulation and enhanced financial stability. Financial markets, including the stock market, are crucial for fostering economic development by facilitating the movement of capital to enterprises and ensuring the effective allocation of resources (Mishkin, 2015). For individual investors, engaging in stocks and other financial instruments can serve as a means to accumulate personal wealth over time. Research indicates that in advanced economies, such as the United States, the stock market significantly contributes to corporate investment, stimulates innovation, and propels GDP growth (Pagano, 1993).

Despite the potential rewards of investing, the financial market is inherently risky, particularly for investors who lack the necessary knowledge to make informed decisions. This situation emphasizes the critical necessity for innovative educational resources and tools aimed at enhancing investment acumen, particularly among younger demographics. The 2024 S&P Global Investment Management Index survey indicates a marked change in investor attitudes, revealing that U.S. investors are becoming increasingly risk-averse. This shift is largely influenced by apprehensions regarding market sustainability in the context of geopolitical tensions, unpredictable political environments, and historically elevated stock valuations (S&P Global, 2024). The survey further reveals a significant decline in the risk appetite index, which has plummeted to -5%—the lowest level recorded since October 2023—down from 28% in May 2024, thereby underscoring the fragility of investor confidence and the potential challenges facing market stability.

The research conducted by Ariel Investments, Charles Schwab, and S&P Global highlights a significant opportunity to develop specialized educational and advisory resources aimed at novice investors. By tackling the deficiencies in financial literacy through robust educational programs, individuals can be empowered to make more informed investment choices, thereby improving their personal financial results and supporting a stable and resilient investment environment. Enhancing financial comprehension, particularly among younger populations, can lead to increased resilience of financial markets in the face of economic fluctuations and uncertainties (Ariel Investments & Charles Schwab, 2022; OECD, 2020; S&P Global, 2024).

**Risk appetite falls to 9-month low**
Risk appetite for next 30 days, net balance (%)



Data accessed July 9, 2024.
The net balance shows the percentage of those risk tolerant minus the percentage of those risk averse. Those reporting a high tolerance or aversion count with double weight.
Source: S&P Global Market Intelligence.
© 2024 S&P Global.

*Figure 1. Risk Appetite Index Trends (2020-2024)*

In Vietnam, the financial market is undergoing significant expansion, driven by a growing interest in investment among the youth demographic. This trend is reflected in the rising engagement with digital financial platforms, which enhance accessibility to investment opportunities. However, akin to patterns observed in the United States, a considerable number of these emerging Vietnamese investors, especially those aged below 35, exhibit a limited understanding of the risks inherent in various financial instruments, including stocks, bonds, and mutual funds (Vietnam Institute of Finance and Investment, 2024). This deficiency in knowledge not only results in poorly informed financial choices but also increases the potential for market instability, as a surge of novice investors may exacerbate volatility in reaction to economic changes.

For Vietnam, implementing structured financial education programs, inspired by successful initiatives in developed countries, could effectively bridge these knowledge gaps, empower young investors to make informed decisions, and foster a robust and sustainable market. By equipping Vietnamese investors with competencies in risk evaluation and long-term financial planning, they will be better positioned to navigate intricate financial environments, thereby promoting individual financial advancement and contributing to the overall stability of the national economy.

## 1.2. Problem Overview

The financial technology (fintech) sector in Vietnam has experienced significant growth, particularly with the emergence of applications designed to assist both novice and experienced investors. Despite the availability of such platforms, new investors in Vietnam have been slow to adopt these technologies. Many individuals, especially those just starting their investment journey, tend to rely on traditional methods or informal advice rather than fintech platforms (Nguyen, 2023). This reluctance stems from several factors, including the complexity of the platforms and a lack of personalized guidance tailored to beginners.

A few prominent platforms in Vietnam, such as SSI iBoard, VnDirect, and FiinPro-X, offer a range of services tailored to the stock market, including real-time stock trading, access to market data, and portfolio management tools. These platforms are designed to provide a comprehensive set of features that cater to both novice and experienced investors. While these platforms provide sophisticated tools for seasoned investors, surveys suggest that new investors perceive them as overly complex and lacking the educational resources needed to build foundational knowledge (Nguyen, 2023). For instance, many of the learning materials and tutorials focus on technical analysis and market predictions, which can be difficult for beginners to grasp without prior financial literacy. Moreover, the platforms often fail to offer personalized investment advice that aligns with individual users' financial goals, risk tolerance, or experience level, making it difficult for new investors to make informed decisions (Tran et al., 2022).

The primary limitation of current fintech platforms lies in their dependence on generalized market insights. Although these platforms utilize artificial intelligence (AI) and big data to analyze stock trends and provide investment recommendations, their suggestions often overlook individual user preferences, such as specific financial goals or risk tolerance (Tran et al., 2022). This lack of personalization, coupled with complex user interfaces, can be a major barrier for novice investors. Many of these users lack the financial knowledge, research skills, and understanding of financial statements required to make informed investment decisions. As a result, they may feel overwhelmed and discouraged by the volume and intricacy of information presented on these platforms. Additionally, these platforms often fall short in providing educational support; few offer accessible, beginner-friendly tools to help new investors understand the complexities of stock market investing. These issues are further compounded by the unique characteristics of the Vietnamese market, which differs considerably from the Western financial systems for which most fintech platforms are designed. Popular U.S.-based platforms like Stash, Acorns, and Robinhood are renowned for their user-friendly interfaces, educational resources, and personalized recommendations, which have made

investing more accessible to individuals with limited financial literacy. However, these platforms primarily target Western stock markets, which have different growth patterns and volatility levels compared to the Vietnamese market (Ariel Investments & Charles Schwab, 2022; S&P Global, 2024). This disparity creates additional challenges for Vietnamese investors, as the insights provided by these platforms may not align with the unique conditions of their domestic market (Vietnam Institute of Finance and Investment, 2023).

To bridge this gap, our project aims to develop a stock price prediction platform tailored specifically for the Vietnamese market. By leveraging artificial intelligence (AI) and machine learning, the platform will provide accurate, real-time predictions while addressing the unique needs of new investors. It will offer personalized investment recommendations that consider users' financial goals, risk tolerance, and experience level, making the process more intuitive and supportive for novices. Through a user-friendly interface, tailored advice, and comprehensive educational tools, this platform will empower Vietnamese investors to make informed financial decisions, facilitating broader participation in the stock market and fostering long-term financial stability (Jones & Smith, 2023; Fleming, 2024).

## 1.3. Motivation of the Project

The motivation behind this project arises from the specific challenges faced by novice investors in Vietnam, who often struggle to navigate the complexities of a fast-growing but volatile stock market. Current platforms lack personalization, offering only generic insights and complex interfaces that are not suitable for beginners, making the investment process even more difficult. Additionally, without personalized financial advice and effective risk management guidance, new investors face heightened risks and barriers to informed decision-making (Nguyen, 2023).

To address these issues, this project aims to create a user-friendly platform tailored to Vietnamese investors. Our solution will leverage machine learning to predict stock trends, provide AI-driven, personalized recommendations, and incorporate educational resources to enhance users' investment knowledge and decision-making capabilities. A key feature of the platform will include a quick research support tool, allowing users to easily search for and analyze relevant information.

The scope of the project focuses primarily on the Vietnamese stock market, with three main components: AI models to predict stock price trends, offering real-time insights; AI solutions for portfolio management, helping users optimize returns while minimizing risks; and an AI-powered chatbot for automated, tailored financial advice. The project

also includes an exploration of the broader implications of AI in finance, with policy recommendations to support AI development in Vietnam's financial sector, aligning with the government's 2030 AI development objectives.

## 1.4. Thesis Outline

This report is organized into several main sections to provide a clear and comprehensive understanding of the topic. It begins with an Introduction that sets the stage, followed by a Literature Review that explores existing research on AI in finance, covering both theoretical insights and empirical evidence. The Methodology section then outlines the methods used for analyzing case studies and gathering relevant data. Next, the Results and Evaluation section presents the findings and assesses their significance. In the Discussion, these results are analyzed in the context of existing literature, with consideration of practical implications. Finally, the Conclusion and Recommendations summarize key takeaways and propose future directions for research and advancements in AI within the financial sector.

## II. Literature Reviews

## 2.1. Investment economics and portfolio management theories

Investment economics has evolved substantially over the past few decades, shaped by the continuous development of portfolio management theories and the integration of artificial intelligence in finance. This review explores these areas, focusing on Expected Utility Theory (EUT), Mean-Variance Portfolio Theory (MVP), and recent advancements in AI for financial applications.

Investment economics is fundamentally concerned with the allocation of resources in a way that maximizes returns while managing risk. One of the foundational theories within this field is the Expected Utility Theory (EUT), developed by Von Neumann and Morgenstern (1944). EUT posits that individuals make investment decisions based on the maximization of their expected utility, a measure of the satisfaction or benefit derived from potential outcomes. According to Barberis (2013), EUT has significantly influenced financial modeling, as it introduces a systematic approach to understanding investor preferences under uncertainty. However, critiques of EUT highlight that it assumes rationality and consistency in decision-making, which may not hold in real-world scenarios due to behavioral biases and psychological factors (Tversky & Kahneman, 1992).

The Mean-Variance Portfolio (MVP) Theory, introduced by Harry Markowitz (1952), builds upon EUT by providing a quantitative framework for portfolio selection that aims

to optimize returns relative to risk. Markowitz's framework uses variance as a proxy for risk and enables investors to construct efficient portfolios, balancing risk and reward by diversifying assets. This model laid the groundwork for modern portfolio theory and is widely used in financial institutions for risk management and asset allocation strategies. While MVP theory is practical, it also assumes that asset returns are normally distributed and investors have constant risk tolerance, which can be limiting in dynamic market conditions (Markowitz, 1959; Elton & Gruber, 1997).

The emergence of Artificial Intelligence (AI) has led to significant progress in the finance sector, fundamentally altering conventional investment strategies. AI technologies facilitate innovative methods for managing extensive data sets, enabling real-time forecasting, and customizing investment recommendations to meet individual requirements. Techniques such as machine learning (ML) and deep learning are now commonly employed to forecast stock prices, enhance trading algorithms, and evaluate financial risks. Importantly, AI systems utilize natural language processing (NLP) to gauge sentiment from news outlets, which has become an essential element in predictive financial models (Baker & Wurgler, 2007). Additionally, AI improves portfolio management through reinforcement learning, optimizing asset distribution, and modeling intricate, non-linear relationships within financial datasets (Li et al., 2019). Among the notable applications of AI in finance are robo-advisory services that offer automated financial planning tailored to the specific preferences of individual investors. These platforms utilize AI algorithms to assess clients' risk profiles and suggest suitable asset allocations, thereby broadening access to investment guidance (Belanche et al., 2019). While AI-driven tools enhance operational efficiency, they also present challenges, including the need for model interpretability, addressing ethical issues, and reducing algorithmic biases (Agrawal, Gans, & Goldfarb, 2018). The dependability of AI systems during uncertain or unprecedented market conditions is a vital area of investigation, as these systems frequently rely on historical data that may not reliably forecast future market behaviors (Bengio, 2019).

The incorporation of Expected Utility Theory (EUT) and Mean Variance Portfolio (MVP) models within traditional investment economics has created a solid theoretical foundation for comprehending investor behavior and risk management. Nevertheless, as financial markets become increasingly complex, AI offers powerful solutions.

## 2.2. Time Series Forecasting

Time-series modeling has historically been a key area of academic research - forming an integral part of applications in topics such as climate modeling (Mudelsee, 2018), biological sciences (Stoffer & Ombao, 2012) and medicine (Topol, 2018), as well as

commercial decision making in retail (Böse et al., 2017) and finance to name a few. While traditional methods have focused on parametric models informed by domain expertise—such as autoregressive (AR), exponential smoothing (Gardner, 1985),(Winters, 1960) or structural time-series models —modern machine learning methods provide a means to learn temporal dynamics in a purely data-driven manner (Ahmed et al., 2010). Deep learning in particular has gained popularity in recent times, inspired by notable achievements in image classification (Ahmed NK, Atiya AF, Gayar NE, El-Shishiny H. 2010), natural language processing (Devlin J, Chang MW, Lee K, Toutanova K. 2019) and reinforcement learning (Silver et al., 2016). By incorporating bespoke architectural assumptions—or inductive biases (Baxter, 2000)—that reflect the nuances of underlying datasets, deep neural networks are able to learn complex data representations (*Representation Learning: A Review and New Perspectives*, 2013), which alleviates the need for manual feature engineering and model design. The availability of open-source backpropagation frameworks has also simplified the network training, allowing for the customization for network components and loss functions. There are several time series models that we can use for forecasting, which can generally be categorized into two different categories:
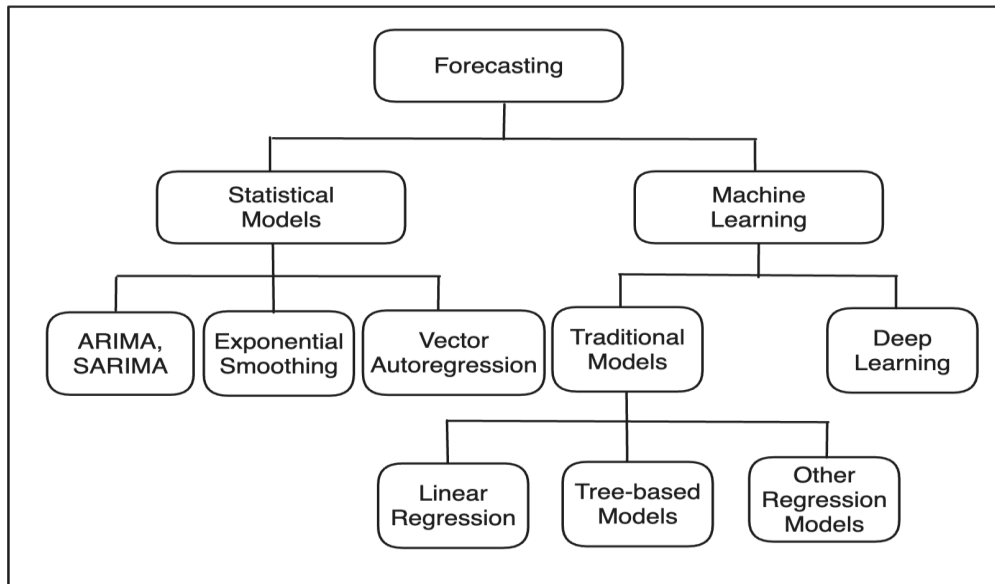


*Figure 2. Classification diagram of forecasting methods*

### 2.2.1. Traditional time series forecasting models

Traditional models offer key advantages in terms of simplicity and interpretability. Models such as ARIMA and SARIMA effectively capture trends and seasonality in data, while Moving Averages and Exponential Smoothing excel at reducing noise, emphasizing core trends (Box, Jenkins, & Reinsel, 2015; Chatfield, 2019). Decision tree models such as XGBoost and LightGBM are also very effective with large data sets, allowing selection of important features to improve predictions (Chen & Guestrin, 2016). With their low computational demands and flexibility, traditional models remain a popular and effective choice for a variety of forecasting applications.

**ARIMA (Auto-Regressive Integrated Moving Average):**
ARIMA is a fundamental statistical model in time series forecasting, featuring three main components: autoregression (AR), integration (I), and moving average (MA). ARIMA is effective for time series data that are linear and stationary. However, it faces challenges when predicting nonlinear or complex time series and requires strong assumptions about the stationarity and properties of the series. Despite its simplicity, ARIMA remains widely used due to its ease of implementation and interpretability.
To understand how the ARIMA model works, there are three terms in its name that you need to understand more clearly:
AutoRegressive - AR(p) is a regression model with lagged values of $y$ up to the p-th lag in the past as predictors. Here, p is the number of lagged observations in the model, $\epsilon$ is the white noise at time t , c is a constant, and $\phi$ are the parameters.

$$\widehat{y}_t = c + \varphi_1 y_{t-1} + \varphi_2 y_{t-2} + \ldots + \varphi_p y_{t-p} + \varepsilon_t$$

**Integrated** I(d) - Differencing is applied d times until the original series becomes stationary. A stationary time series is one whose properties do not depend on the time at which the series is observed.

$$By_t = y_{t-1} \text{ where B is called a backshift operator}$$

Thus, a first order difference is written as:

$$y'_t = y_t - y_{t-1} = (1 - B)y_t$$

In general, *a d th*-order difference can be written as:

$$y'_t = (1 - B)^d y_t$$

**Moving Average MA(q)** - The moving average model uses a regression-like approach on past forecast errors. Here, ε represents the white noise at time t, c is a constant, and θs are the parameters.

$$\hat{y}_t = c + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + ... + \theta_q \varepsilon_{t-q}$$

Combining all three types of models above will produce the ARIMA(p,d,q) model.

$$\hat{y}'_t = c + \phi_1 y'_{t-1} + \phi_2 y'_{t-2} + ... + \phi_p y'_{t-p} + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + ... + \theta_q \varepsilon_{t-q} + \varepsilon_t$$

**XGBoost (Extreme Gradient Boosting):**
XGBoost is a machine learning algorithm based on the gradient boosting technique, known for its computational efficiency and high predictive power. XGBoost builds decision trees incrementally to minimize prediction errors, making it suitable for both linear and nonlinear data, particularly in large and complex datasets. A drawback of XGBoost is the extensive hyperparameter tuning required for optimal performance, which can demand considerable computational resources when applied to very large datasets.
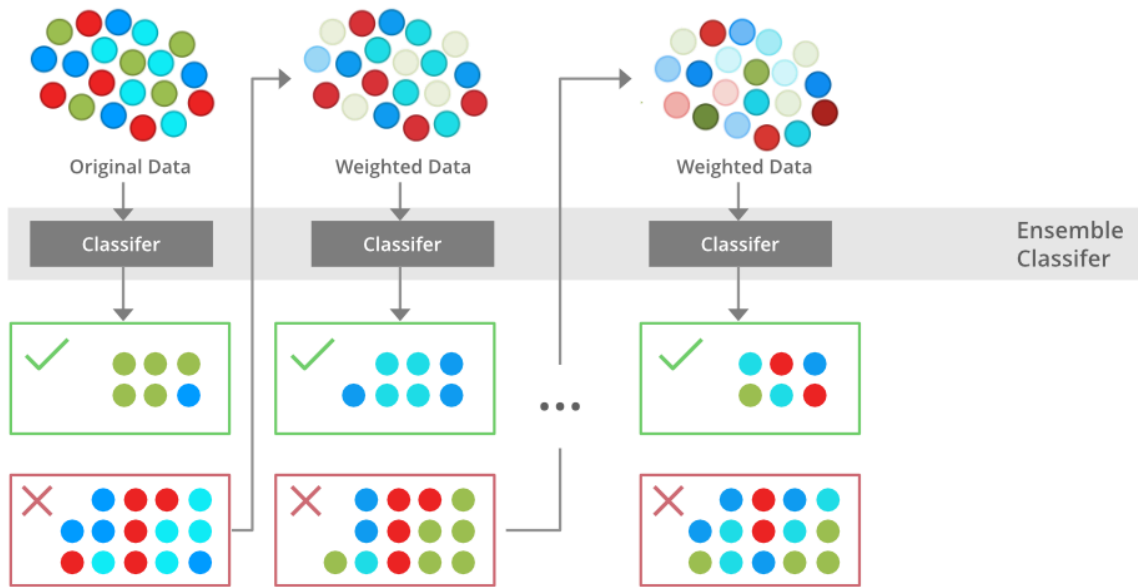


*Figure 3. Explaining how the XGBoost model works*

**LightGBM (Light Gradient Boosting Machine):**

LightGBM, as a framework to implement GBDT algorithm, combines GOSS algorithm and EFB algorithm and is used in data sets with large sample data and high dimensions (Yang, Y., Wu, Y., Wang, P., & Jiali, X. (2021)). Its optimization features include Leaf-Wise based decision tree growth strategy, optimal segmentation of category eigenvalues, feature parallelism and data parallelism.It not only reduces the communication overhead and time complexity between data, but also ensures the rationality of data



*Figure 4.* Illustration of LightGBM model

### 2.2.2. Deep learning models

Deep learning models like RNNs, LSTMs, and Transformer-based models excel in capturing complex and long-term dependencies in time-series data. RNNs were pioneering models for sequential data, while LSTMs improved them by addressing the vanishing gradient issue, allowing for longer memory in forecasting (Goodfellow, Bengio, & Courville, 2016; Hochreiter & Schmidhuber, 1997). Transformers, with their self-attention mechanism, handle dependencies across long timeframes without sequential processing, making them highly efficient and accurate on large datasets (Vaswani et al., 2017).

**LSTM (Long Short-Term Memory):**

LSTM is a type of recurrent neural network (RNN) specifically designed to handle long-term dependencies in sequential data. With its memory cell structure, LSTM can retain information over extended time intervals, making it highly effective in forecasting non-stationary and nonlinear time series. This capability has made LSTM one of the top choices in forecasting applications across finance, manufacturing, and other fields. However, a significant limitation of LSTM is its requirement for substantial computational resources and long training times.



***Figure 5.*** *LSTM model illustration*

## 2.3. LLM Pre-trained model:

In recent years, Natural Language Processing (NLP) has become a foundational technology for developing sophisticated chatbot systems, largely due to advancements in pre-trained models like Large Language Models (LLMs), Optical Character Recognition (OCR), and embedding models. These pre-trained models significantly enhance chatbots' ability to understand and generate text that closely mimics human communication. Prominent LLMs such as GPT-3, BERT, and T5 have been widely applied to tasks like generating responses, semantic analysis, and dialogue management (Brown et al., 2020; Devlin et al., 2019). Embedding models further improve a chatbot's ability to understand context and carry out tasks like information retrieval, sentiment analysis, and recommendations, leading to richer user interactions (Reimers & Gurevych, 2019).

A notable application in this domain is the development of a Vietnamese Question Answering System, which transitioned from utilizing multilingual BERT models to a

specialized monolingual BERT model. This shift enhances the model's performance in understanding and processing Vietnamese language nuances, thereby improving the accuracy of responses in a localized context (Pham et al., 2021). Despite these advancements, the use of LLMs in chatbot applications still faces challenges such as response hallucinations and inaccuracies. These limitations have led to the development of Retrieval-Augmented Generation (RAG) models to address such concerns (Lewis et al., 2020).

### 2.3.1. Q&A SQL and Q&A RAG

One prominent approach for improving the accuracy of chatbot systems interacting with databases is SQL-based query generation. SQLNet by Li et al. (2018) introduced a model that enables LLMs to generate structured SQL queries from natural language questions without relying on reinforcement learning. Their framework demonstrated that LLMs could be trained to accurately interpret user intent and translate it into executable SQL statements, especially useful in cases where chatbots are required to interact with large relational databases. This approach underscores the potential for LLM-based systems to be adapted for data-driven environments, though challenges in real-time optimization persist (Chen et al., 2021).



**Figure 6.** *Building Q&A systems of SQL databases requires executing model-generated SQL queries.*

In addition to SQL-based query generation, Retrieve-and-Generate (RAG) models have also contributed significantly to the development of knowledge-intensive chatbot systems. RAG frameworks, as proposed by Lewis et al. (2020), are designed to combine the benefits of information retrieval with language generation, enabling chatbots to handle both structured and unstructured data sources. By retrieving relevant information from external databases and generating a coherent response, RAG models address the limitations of purely generative models that may otherwise hallucinate facts. This dual

approach has proven particularly effective in handling complex user queries where both precise data retrieval and language fluidity are essential (Zhang & Shi, 2023).



**Figure 7.** *QA RAG architecture and workflow. Source: (Mansurova et al., 2024c)*

Further research has explored the integration of SQL Agents with RAG models and LLMs to enhance the precision and contextual relevance of chatbot responses. Chen et al. (2021) discuss a multi-agent framework where a SQL Agent acts as an intermediary, optimizing SQL queries generated by LLMs and refining them based on user feedback. This integration is essential for applications requiring high accuracy in data retrieval and real-time response generation, especially when chatbots interact with databases with dynamically updated content. Their findings indicate that SQL Agents can reduce processing time while maintaining high query accuracy, thus improving the overall efficiency of chatbot interactions.

### 2.3.2. Agent-graph and Multi agent:

Agent-graph and multi-agent systems represent a growing area of research in chatbot development due to their capability to manage complex and distributed tasks. In essence, multi-agent systems (MAS) involve a network of agents working collaboratively to solve problems that are beyond the capacity of a single agent. Each agent within the system is designed to perform specific functions autonomously while communicating and cooperating with other agents to achieve a collective goal (Ferber, 1999). Agent-graph

models are particularly advantageous when structuring and visualizing the interactions and workflows between agents. In this model, agents are represented as nodes, and their interactions are depicted as edges within a graph. Such a representation allows for a clearer understanding of how different agents coordinate tasks, exchange information, and respond to user inputs. For instance, in a chatbot application designed for stock market advisory, different agents may handle specific functions - one may focus on real-time stock data retrieval, another on processing user inquiries, while others may analyze trends and generate personalized investment advice (Wooldridge, 2009). This multi-layered interaction framework enhances the chatbot's ability to respond quickly and accurately to a wide range of user queries.



***Figure 8.*** *Multi-agent system architecture*

The advantages of using multi-agent systems in chatbots extend beyond task division and autonomy. These systems promote scalability, as more agents can be added to the network without compromising overall efficiency. Additionally, agent-graph models ensure that failures in individual agents do not affect the system's entire functionality. By distributing tasks among specialized agents, MAS-based chatbots can balance workloads and reduce bottlenecks during peak usage, an essential factor in high-demand services like financial advisory (Jennings & Wooldridge, 1995).

However, challenges remain, particularly regarding coordination and communication between agents, which is crucial for maintaining an effective multi-agent system.

Solutions often involve adopting protocols for agent cooperation, such as the Contract Net Protocol (Smith, 1980), or employing decentralized approaches to reduce dependency on central control. Moreover, the development of trust and security measures within agent networks is necessary to ensure that users receive accurate and reliable information, especially in sensitive domains like finance (Calvaresi et al., 2017). As MAS continues to evolve, it is expected to drive further improvements in the robustness and adaptability of chatbot systems.

In conclusion, the literature highlights significant advancements in chatbot development through the integration of LLMs, RAG models, Langsmith and the handling of structured data with Agent - graph, Multi agents. At the same time, it underscores the importance of robust evaluation frameworks to ensure these systems meet the performance expectations of users and deliver reliable, accurate interactions.

## 2.4. Evaluation Metrics

### 2.4.1. Time-series prediction

Mean Squared Error (MSE) is perhaps the most commonly used metric for regression problems. Essentially, it finds the average squared error between predicted and actual values. MSE is a measure of the quality of an estimator—it is always non-negative, and the closer the values are to 0, the better.(Piotrowski, P., Rutyna, I., Baczyński, D., & Kopyt, M. (2022))

$$MSE = \frac{\sum (y_i - \hat{y}_i)^2}{n}$$

Where n is the number of data points, $y_i$ is the observed value, and $\hat{y}_i$ is the predicted value.

In regression analysis, plotting is a more natural way to see the overall trend of the data. Simply put, MSE tells you how close the regression line is to a set of points. It does this by measuring the distance from the points to the regression line (these distances are the "errors") and squaring them. Squaring is important to avoid complications with negative signs. It also places more weight on larger differences.

To minimize MSE, the model can become more accurate, meaning the model is closer to the actual data. An example of linear regression using this method is the least squares

method, which evaluates the fit of a linear regression model to a two-variable dataset, although its limitations are related to the known distribution of the data. The lower the MSE, the better the forecast.

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \widehat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2}$$

The second metric we use in the project, $R^2$ metric - also known as the coefficient of determination — is an important measure of how well a model fits the data, as well as its ability to predict future outcomes. In simple terms, $R^2$ indicates how much of the variation in your data can be explained by the model. The closer the $R^2$ value is to 1, the better the model fits the data, and the stronger the model's predictive ability.

The $R^2$ metric is calculated using the formula:

In the field of stock price prediction, $R^2$ helps evaluate how well a model explains stock price fluctuations using historical data. A high $R^2$ value means the model is better able to capture market trends and accurately reflect price changes in past data, indicating that the model can explain a large portion of the variation in stock prices.

Another metric we use to evaluate the model is the Mean Absolute Error (MAE), which is a commonly used metric to evaluate the accuracy of regression models by measuring the average magnitude of the errors between the predicted and the actual values, without considering their direction. Unlike the Mean Squared Error (MSE), which is the square of the errors, the MAE calculates the absolute difference between the predicted values $(\widehat{y}_i)$ and observed values $(y_i)$, making it a straightforward representation of error in the same unit as the data.

**Formula:**

$$MAE = \frac{1}{n}\sum_{i-1}^{n}\left|y_i - \widehat{y}_i\right|$$

where:

- *n is the total number of data points,*

- $y_i$ is the observed value, and
- $\widehat{y}_i$ is the predicted value.

In stock price prediction, MAE can help analysts understand the typical prediction error. A model with a low MAE on historical data will have a higher reliability in predicting future trends. While MAE may not capture extreme deviations as strongly as MSE, it provides an accurate estimate of general forecasting accuracy, which can be especially useful in minimizing consistent prediction errors in the financial domain.

### 2.4.2. About NLP Pre-train models

In assessing the efficacy of chatbot systems, a variety of metrics are utilized to evaluate both linguistic and operational performance. Notably, BLEU (Bilingual Evaluation Understudy) and ROUGE (Recall-Oriented Understudy for Gisting Evaluation) scores are frequently employed to gauge the quality of the text generated by chatbots in relation to human responses, thereby offering insights into the language generation capabilities of these systems (Papineni et al., 2002; Lin, 2004).

$$BLEU = P. exp( \sum_{n=1}^{N} w_n \, log \, p_n )$$

- $w_n$: positive weights for n-grams, summing to 1
- $p_n$: n-gram precision

Furthermore, metrics tailored to specific tasks, such as success rate, user satisfaction, and real-time performance, are essential for determining a chatbot's practical effectiveness, especially in high-stakes contexts like financial decision-making (Zhou et al., 2020). The Goal-Question-Metric (GQM) framework is also commonly applied to evaluate multi-agent chatbots, providing a systematic method for measuring the system's readiness, adaptability, and overall effectiveness across various application domains (Kitchenham et al., 2002; Galster et al., 2014).

LangSmith is a noteworthy platform designed specifically for tracking and assessing the effectiveness of chatbots in real-time applications. LangSmith provides a comprehensive suite of tools that automate performance monitoring, making it easier to analyze key metrics such as response accuracy, latency, and user satisfaction over time. This automation helps developers maintain high service quality, even as chatbots handle

dynamic and large-scale user interactions (Smith et al., 2022). One of LangSmith's most significant features is its ability to evaluate conversational flows and detect potential issues, including misinterpretations or delays in chatbot responses. By monitoring these aspects, LangSmith allows developers to fine-tune NLP algorithms, thereby improving both the accuracy and contextual relevance of chatbot replies (Patel & Yang, 2021). Additionally, LangSmith integrates well with popular NLP pre-trained models, such as GPT and BERT, supporting a seamless process to track model efficacy across different user intents and sentiment types. LangSmith's metrics also align with industry-standard evaluation parameters, making it a valuable resource for projects aiming to assess and optimize chatbot performance effectively.

## III. Methodology

### 3.1. Datasets

In this project, we utilized two primary data streams: data from Application Programming Interfaces (APIs) and data from the user input and external documents for the investment advisory chatbot. Specifically, the data provided by APIs enabled us to retrieve real-time stock market data and information about listed companies. This data plays a pivotal role in the financial forecasting models and supports our investment advisory chatbot.

For real-time access to Vietnamese stock market data, we integrated the **VNStock** API, an open-source Python library hosted on GitHub, developed by **Thinh Vu** (n.d.). VNStock offers seamless access to data from the three major stock exchanges in Vietnam: Ho Chi Minh Stock Exchange (HOSE), Hanoi Stock Exchange (HNX), and Unlisted Public Companies Market (UPCoM). This API provides essential stock-related information such as stock prices, trading volumes, and various key financial indicators that are crucial for both financial forecasting and the AI chatbot.

VNStock API is highly efficient in handling large amounts of data, enabling us to automate the process of data retrieval from multiple exchanges and ensure our models are constantly updated with the latest market information. This integration is especially valuable for our project, as it reduces the reliance on manual data collection, ensuring data consistency, and allowing the system to make real-time decisions based on dynamic market changes (Vu, n.d.).

The core data that we retrieve from the VNStock API includes multiple financial indicators that are critical for building and training AI models capable of predicting future stock movements. The following are the key data points:

| | VNM | | | | | VCB | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Open | Close | High | Low | Volume | Open | Close | High | Low | Volume |
| Count | 52539.0 | 55270.0 | 52539.0 | 52539.0 | 5.527000e+04 | 52539.0 | 52539.0 | 52539.0 | 52539.0 | 5.253900e+04 |
| Mean | 68.725187 | 68.724610 | 68.765920 | 68.685196 | 1.674861e+04 | 89.75164 | 89.751802 | 89.792272 | 89.711955 | 6.840062e+03 |
| Std | 2.701632 | 2.701648 | 2.703894 | 2.697624 | 3.780903e+04 | 3.48850 | 3.488046 | 3.491359 | 3.484637 | 2.107602e+04 |
| Min | 63.1 | 63.1 | 63.2 | 63.0 | 1.000000e+02 | 80.1 | 80.1 | 80.2 | 80.1 | 1.000000e+02 |
| 25% | 66.9 | 66.9 | 66.9 | 66.9 | 2.900000e+03 | 87.60000 | 87.6 | 87.7 | 87.6 | 1.100000e+03 |
| 50% | 67.9 | 67.9 | 68.0 | 67.9 | 7.500000e+03 | 89.80000 | 89.8 | 89.8 | 89.8 | 3.100000e+03 |
| 75% | 70.4 | 70.4 | 70.4 | 70.3 | 1.750000e+04 | 92.00000 | 92.0 | 92.1 | 92.0 | 7.300000e+03 |
| Max | 76.1 | 76.1 | 76.2 | 76.0 | 2.744300e+06 | 100.50000 | 100.5 | 100.5 | 100.3 | 2.202800e+06 |

**Table 1.** *Statistics on VNM and VCB stock prices. Source: Made by the author's group*

By utilizing these indicators, our predictive models leverage historical and current data to make informed predictions about future stock movements. The integration of such indicators ensures that our AI models are trained on high-quality, relevant data, thus improving their accuracy and reliability in providing financial advice and forecasts.

In addition to stock market data, the project collects **company-specific information** such as news articles, major company events, and financial reports. This unstructured data is processed using Natural Language Processing (NLP) techniques, which help our models assess market sentiment and the potential impact of company-specific events on stock prices (Tetlock, 2007). The integration of company-specific data allows the system to provide a more comprehensive analysis of individual stocks, incorporating both quantitative data (stock prices, volumes) and qualitative data (company news, sentiment).

To complement stock market data, we incorporated user-provided data and external document sources, which help the chatbot provide more contextualized and comprehensive advice. This data includes both structured and unstructured formats:

**PDF Files and Scanned Documents**: These files are processed using Optical Character Recognition (OCR) technology, which converts them into searchable text. OCR enables the chatbot to retrieve relevant information from various document formats, making it accessible for user queries (Smith, 2021).

**Structured Files (CSV and XLS)**: Structured files, such as CSV and XLS formats, provide organized datasets that the chatbot can quickly reference for specific queries. This predictable structure allows for fast and efficient data retrieval (Johnson, 2020).

**SQL Databases**: We utilize PostgreSQL as the primary relational database management system for storing and managing structured data. PostgreSQL's robust functionality ensures efficient data organization, supporting rapid data retrieval and manipulation as needed for chatbot interactions (Doe, 2023).

**Web Sources**: The chatbot also retrieves real-time information from the web to respond dynamically to user inquiries, ensuring that users receive accurate and current information (Brown & Green, 2022).

Through a combination of real-time API data and structured/unstructured data from user documents, we establish a robust, versatile knowledge base. This allows our system to offer more insightful predictions and advice, catering to a diverse range of user inquiries

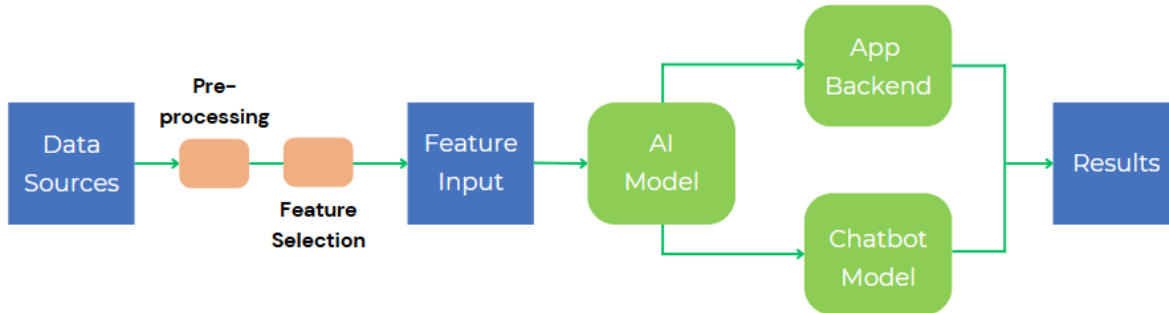and supporting more accurate financial decision-making.

## 3.2. Proposed Method



***Figure 9.*** *Overall Proposed Method. Source: Made by the author's group*

The data processing phase is crucial for transforming raw data into structured formats that optimize the performance of our financial forecasting models and investment advisory chatbot. This phase begins with data collection, where we integrate data from both the internal management system and the VNStock API. The internal management system provides historical data, including client investment histories and feedback, allowing the chatbot to generate personalized investment insights. Meanwhile, the VNStock API—an open-source Python library developed by Thinh Vu—delivers real-time stock market data from Vietnam's main stock exchanges (Vu, n.d.). Through this API, we retrieve critical financial indicators such as open, close, high, and low prices, along with trading volume and company-specific updates. The combination of structured financial data and unstructured company information ensures our AI models are equipped with the most recent market information, enhancing the accuracy of the chatbot's recommendations.

Once collected, data is organized in a structured repository, ensuring efficiency and reliability in data access. Organized through data warehousing, our repository holds financial indicators in standardized tables, while company-specific information is stored in unstructured formats, facilitating easy retrieval and contextual analysis (Nguyen, 2020). Data is also segmented by stock symbols, temporal intervals, and indices, while a real-time data pipeline supports continuous data updates, reinforcing data accuracy and timeliness.

In the data preprocessing stage, we prepare the data for analysis by employing cleaning, normalization, and feature engineering techniques. Data cleaning addresses missing

values and outliers, with statistical imputation used to replace missing entries. Normalization and scaling methods, such as Min-Max scaling and Z-score normalization, are applied to adjust for wide variations in stock prices and volumes, ensuring compatibility across algorithms. Additionally, feature engineering enhances data inputs with new metrics: technical indicators such as Simple Moving Averages (SMA), Exponential Moving Averages (EMA), and Bollinger Bands enrich our predictive capabilities, while sentiment scores generated from company news offer insights into investor sentiment. Once the data preprocessing is complete, we will integrate these steps into an automated data processing workflow. This workflow will be triggered each time new data is retrieved from the VNStock API, ensuring that the data is continuously updated and preprocessed, thus maintaining the efficiency and relevance of our financial forecasting models and the investment advisory chatbot.

The training phase encompasses two primary AI models: a time series model for stock price prediction and a Natural Language Processing (NLP) retrained model for the chatbot system. The time series model, utilizing a neural network architecture based on Long Short-Term Memory (LSTM), is specifically designed to forecast stock prices. This architecture features two LSTM layers supplemented with Dropout layers for regularization, culminating in a Dense layer that generates the final output for stock price predictions. The predictions made by this model are integrated into the application system, facilitating real-time price forecasting for users. Each component of LSTM is described in detail below:
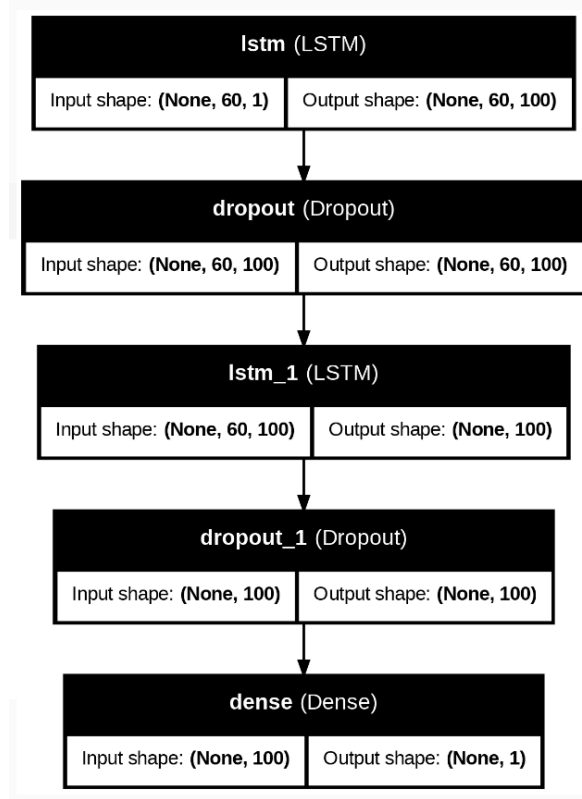
***Figure 10****. LSTM Model Architecture*

Conversely, the NLP model is integrated into the chatbot system, which employs a multi-agent architecture based on the Follower-Leader framework. In this structure, a central "Leader" agent coordinates the actions of multiple "Follower" agents, each assigned to manage specific facets of user queries. The Leader acts as the orchestrator, facilitating the flow of information and making informed decisions based on insights gathered from the Followers. Each Follower agent focuses on processing a distinct aspect of the user input or analyzing different dimensions of a query, thereby enhancing the system's ability to provide tailored responses.

The chatbot system boasts several advanced features designed to optimize user interactions. First, the Q&A with Database Access feature enables the chatbot to retrieve answers by querying a structured database, ensuring accurate and relevant information delivery. Second, the Function Calling for Calculations feature allows the chatbot to execute specific functions or calculations as necessitated by user queries, thereby providing immediate support for user requests. Third, the Retrieval-Augmented Generation (RAG) with Documents functionality enhances response accuracy by employing document retrieval techniques to incorporate pertinent information into the answers provided. Fourth, the Image-Based Information Extraction capability utilizes

optical character recognition (OCR) to extract data from images, broadening the types of queries the chatbot can effectively address. Finally, the Web Search Capability allows the chatbot to pull current information from the web, enriching its responses with up-to-date data.
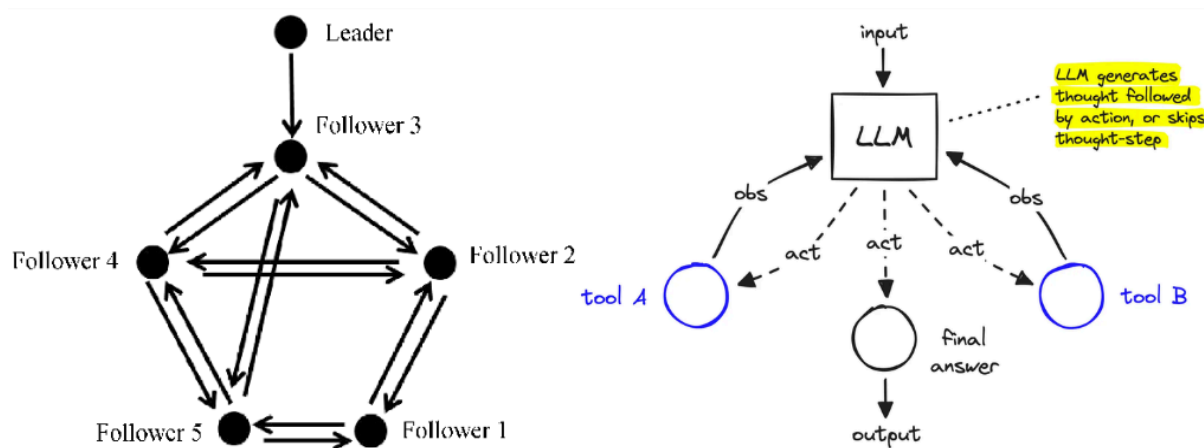


***Figure 11.*** *Follower-Leader framework for Chatbot Architecture.*

Together, these models and features create a robust framework for stock price prediction and user interaction, allowing for precise forecasting and comprehensive query management within the application and chatbot systems. The synergy between the time series model and the NLP chatbot ensures a holistic approach to delivering financial insights and enhancing user experience.

To ensure that the chatbot's responses are quick and accurate, we propose developing a chatbot system architecture instead of solely training a large language model on a pure dataset. This approach will enable the chatbot to answer questions unrelated to the main topic and engage with users in a natural manner while maintaining the clarity of the conversational context. Furthermore, with the architecture we propose, the process of training and improving the model will become easier than ever, as this process will be fully automated.

## 3.3. Experiment Setup

### 3.3.1. LSTM stock prediction

**Runtime environment and libraries**

The experiment was conducted on Google Colab, utilizing an NVIDIA Tesla K80 with 12GB VRAM to enhance LSTM (Long Short-Term Memory) model performance and reduce training time, crucial for computation-intensive time series models (Pascanu et al., 2013). The specific hardware setup includes an Intel Core i7-12500H CPU for efficient data processing, an NVIDIA RTX 3060 with 10GB VRAM for parallel computation and deep learning tasks, and 32GB DDR4 RAM for handling large datasets. This advanced hardware configuration ensures efficient and stable stock price predictions. The model was developed in Python 3.9, widely used in data science and machine learning, with robust support from libraries such as yFinance for historical stock data collection (Tanha et al., 2021), Pandas and NumPy for data processing and time series analysis, Keras (via TensorFlow) for building and training the LSTM model (Chollet et al., 2015), and Matplotlib for data visualization to facilitate comparison between actual and predicted values.

**Data processing strategy**

The data preprocessing strategy was implemented to optimize the LSTM model for stock price prediction. Data was collected from the yFinance API, with the primary feature selected being the closing price (Close), which reflects end-of-day stock fluctuations and helps reduce short-term noise (Shen et al., 2020). To meet LSTM input requirements, data was transformed into fixed-length sequences and normalized using Min-Max Scaling, which scales values to the [0,1] range and enhances training efficiency (Ioffe & Szegedy, 2015). The time series was divided into 60-day windows, with each window comprising 59 days as input and one day as the label, enabling the model to make predictions based on short-term trends (Sethi & Mittal, 2018). Finally, the data was split 80-20 for training and testing to assess the model's performance on unseen, real-world data (Goodfellow et al., 2016).
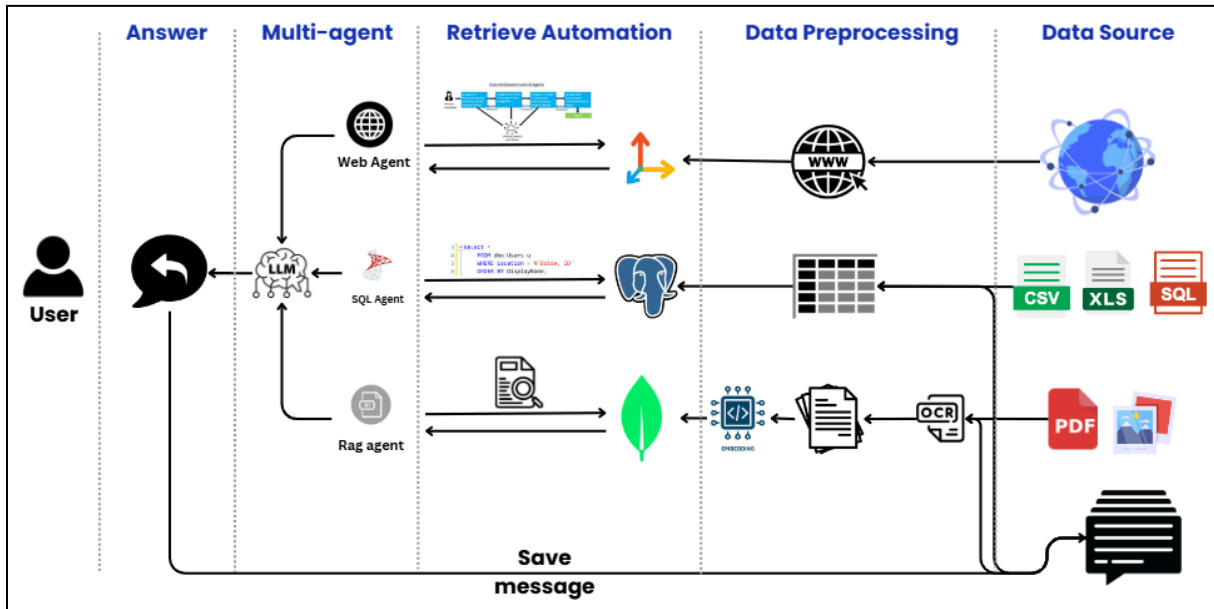
### 3.3.2. Chatbot system architecture



***Figure 12.*** *The process of chatbot system architecture. Source: Made by the author's group*

In terms of system architecture, we have developed an agent graph architecture where agents perform specialized processing tasks across three different datasets, with Vietnam-SBERT serving as the embedding model and the central model. Furthermore, to enable the model to automatically and efficiently retrieve data, we will store relational data in PostgreSQL. Notably, text data, such as textual passages and conversation histories, will be stored in MongoDB. Utilizing different database management systems for various data types can enhance retrieval speed and ensure that the data returned to the chatbot is both accurate and prompt.

Additionally, to allow the chatbot to access and process publicly available data on social media, a dedicated agent will automatically retrieve information using an open-source software called Tavily. Upon receiving an input question from the user, an automated process will be initiated. This process first involves gathering relevant data from various sources, including PDF files, image files, and web data. Subsequently, a data processing workflow will be executed for each type of data, stored using different database management systems. Following this, the agent graph chatbot model will operate by retrieving pertinent information from the relevant databases based on the user's inquiry, thereby facilitating rapid and accurate responses.
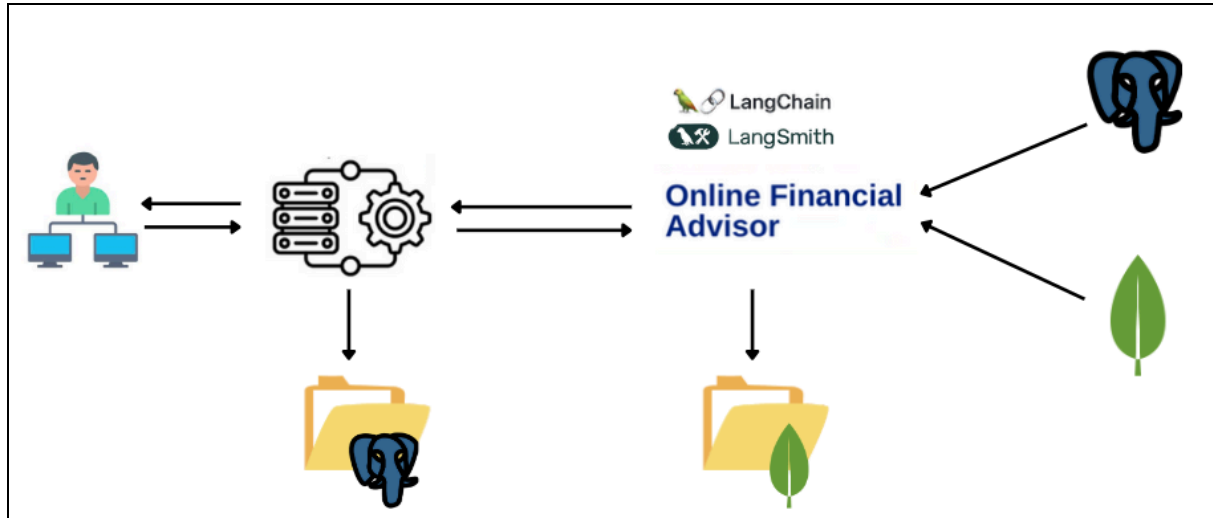
***Figure 13.*** *The architecture of interaction between the chatbot, the application, and the database. Source: Made by the author's group*

To ensure that the architecture operates safely and efficiently, the chatbot system will be specifically developed with a dedicated backend. This design will prevent the backend from becoming overloaded when processing user response tasks. After each interaction with the user, the chatbot's response history will be automatically downloaded to PostgreSQL to facilitate performance processing and evaluation. Additionally, the historical messages will be automatically stored in MongoDB, allowing the chatbot to retrieve previous messages as needed. To ensure that the chatbot cannot access any private data from the application, the databases in each database management system will be separate from the database containing user information.

Based on the aforementioned designs and descriptions, this project will be built on the LangChain platform, one of the most popular frameworks for developing, evaluating, and deploying chatbots today. The system will be evaluated and monitored in detail using LangSmith.

### 3.3.3. System architecture

The frontend is primarily constructed using HTML, CSS, and JavaScript, incorporating libraries such as React and Vue.js. React is chosen for its capability to create reusable UI components, enhancing maintainability and enabling efficient state management (Bergström et al., 2019). Complementing React, Vue.js offers a progressive framework for developing complex user interfaces with minimal setup and high flexibility. jQuery is employed for DOM manipulation and AJAX requests, which enhances user experience by allowing asynchronous data fetching and smooth UI interactions. To ensure a

responsive design, Bootstrap is utilized, enabling the application to adapt effectively across a range of devices and screen sizes.

On the backend, Django serves as the primary framework due to its "batteries-included" philosophy, which supports rapid development and the deployment of secure web applications (Mackenzie, 2021). Django's modular architecture allows for the efficient building of scalable applications. Middleware features enhance security and performance by enabling pre-processing of requests, user authentication, and logging, thus ensuring a robust backend environment. The Django REST Framework (DRF) is also employed to simplify the creation of RESTful APIs, facilitating smooth communication between the frontend and backend.

Data storage and management are handled by PostgreSQL, a robust relational SQL database. Its advanced features, such as ACID compliance and strong data integrity measures, are vital for managing sensitive information like user profiles and financial transactions. The database schema is meticulously designed to optimize data retrieval while maintaining referential integrity across tables.



*Figure 14.* *System Architecture Overview. Source: Made by the author's group*

Redis serves as an in-memory data store, designed to cache frequently requested data, which minimizes the load on the database and enhances user response times (Kim et al., 2022). By integrating Redis with Celery, the system efficiently manages background tasks, particularly for asynchronous data fetching and scheduled updates. This combination not only optimizes overall system performance but also facilitates the handling of multiple concurrent tasks without compromising user experience. To meet

real-time data requirements, such as displaying live stock prices, WebSockets are implemented to create persistent connections between the client and server. This allows for instantaneous data updates, enabling users to receive timely and accurate information. The integration of Celery with Redis supports background job processing, allowing the system to handle high traffic loads without degrading application performance. For example, Celery can be utilized to schedule and execute tasks such as fetching financial data from third-party APIs or conducting periodic data analysis (Johnson et al., 2021).



*Figure 15. Backend - Vnstock Architecture. Source: Made by the author's group*

Furthermore, the backend includes a module that scrapes data from Vnstock, leveraging Redis and Celery for efficient management of these data-scraping tasks. This setup not only automates the collection and updating of information but also ensures that the data remains fresh and accurate. Finally, the use of Docker in the deployment process ensures consistency and scalability of the application, while WebSockets facilitate real-time communication, allowing users to quickly access the information they need. The establishment of supporting frameworks for connectivity helps prevent backend overload from continuously updating data from Vnstock. This approach can enable the application to operate efficiently, handling user requests promptly without any delays.

## IV.    Experiment & Result

### 4.1.    Model Performance Evaluation

Stock price prediction is a complex and challenging task due to the inherently volatile and nonlinear nature of financial markets. Accurate prediction requires robust models that are capable of capturing complex patterns and trends from historical data. In

this section, we analyze and evaluate the performance of four models—ARIMA, XGBoost, LightGBM, and LSTM—in stock price prediction, using Mean Squared Error (MSE), Mean Absolute Error (MAE), and $R^2$ metrics (Table 2) to evaluate their accuracy and efficiency with different optimizers.

| Model | Optimizer | MSE | MAE | $R^2$ |
|---|---|---|---|---|
| Arima | - | 24.600 | 4.10 | 0.75 |
| XGBoost | Adam | 19.980 | 3.75 | 0.79 |
|  | SGD | 20.45 | 3.80 | 0.78 |
| LightGBM | Adam | 20.500 | 3.85 | 0.78 |
|  | SGD | 21.200 | 3.90 | 0.77 |
| LSTM | Adam | 14.300 | 3.10 | 0.87 |
|  | SGD | 17.200 | 3.40 | 0.82 |

***Table 2.** Model performance evaluation. Source: Made by the author's group*

Experimental results show that the LSTM model using the Adam optimizer achieves the highest performance in the stock price prediction problem, with the lowest MSE (14.3), the lowest MAE (3.10), and the highest R² (0.87). With the ability to memorize periodic patterns and model nonlinear trends, LSTM excels in forecasting complex time series, especially with stock price data that is highly volatile. The Adam optimizer helps the model converge faster and more accurately, significantly improving the reliability of the prediction. Meanwhile, XGBoost and LightGBM outperformed ARIMA with R² of 0.79 and 0.78, respectively, but were still inferior to LSTM due to its inability to store time series states, which limits its ability to capture nonlinear relationships. ARIMA performed the worst with an R² of only 0.75, reflecting its limitations in modeling nonlinear trends. Overall, LSTM - a deep learning model - is the most suitable model for accurately forecasting stock prices, outperforming traditional models such as ARIMA and boosting models such as XGBoost and LightGBM.

## 4.2. Testing and Demo WEB

Firstly, this performance evaluation table needs to be adjusted with more detailed parameters to better reflect the actual system requirements. Following the integration of the Celery architecture for automated scheduling, there has been a notable enhancement in the backend's performance when managing requests from VNstock. This setup not only enhances response times but also reduces system strain during times of elevated concurrent requests. By utilizing Celery, the system can process data in real-time, ensuring accurate and prompt stock price predictions.

| Metrics | Desktop | Mobile | Tablet | Target Value | Description |
|---|---|---|---|---|---|
| Page load time (s) | 2.5 | 3.2 | 2.8 | < 3.0 | Average time taken to load the main page fully |
| API Response Time (ms) | 180 | 220 | 210 | < 200 | Time for the chatbot API to respond |
| Prediction Accuracy (%) | 92.1 | 91.7 | 92.0 | > 90 | Average accuracy of stock prediction model |
| Chatbot Response Rate (%) | 98.5 | 97.8 | 98.2 | > 95 | Success rate of chatbot responding correctly |
| User Satisfaction Score | 4.6 | 4.5 | 4.5 | > 4.0 | Average satisfaction score from user feedback |
| Error Rate (%) | 0.8 | 1.1 | 1.0 | < 1 | Percentage of errors encountered by users |
| Session Duration (mins) | 5.2 | 4.9 | 5.1 | - | Average time users stay active on the app |
| Retention Rate (%) | 80.5 | 78.2 | 79.6 | > 75 | Percentage of users who return to the web app |

*Table 3. Web application performance and user experience evaluation metrics. Source: Made by the author's group*

Allocating web resources and optimizing the architecture in this way maintains stable performance across desktops, mobile devices, and tablets. The evaluation table below outlines key performance metrics, such as page load time, prediction accuracy, and user satisfaction rate. These metrics indicate that the application operates quickly and meets user needs effectively. For example, page load times across devices range only from 2.5 to 3.2 seconds, meeting the target of under 3 seconds. The prediction model accuracy exceeds 90%, while the average user satisfaction score is 4.5 out of 5, demonstrating a high level of user satisfaction.

These experimental results confirm that the application is not only fast but also accurate and reliable, enabling users to access information and forecasts effectively and conveniently. The metrics show that the system has met its initial goals for performance and user experience, paving the way for potential upgrades and feature expansion in the future.

## V.    Discussion

The investment advisory system created in this project signifies a considerable leap forward in leveraging artificial intelligence and machine learning to aid novice investors in navigating the intricacies of stock market investments. By incorporating sophisticated predictive models such as ARIMA, LSTM, XGBoost, and LGBM, the system seeks to improve the precision of stock price predictions. Among these models, LSTM is particularly distinguished for its effectiveness with time series data, facilitating a deeper comprehension of price fluctuations based on historical patterns.

A key component of this system is the AI-powered chatbot, which offers users real-time financial guidance. This interactive feature is intended to enrich the user experience by making financial insights more approachable and engaging for individuals who may find traditional investment strategies daunting. The chatbot streamlines the process of acquiring investment advice, empowering novice investors to make more informed decisions with greater confidence.

However, the project faces several significant limitations that need to be addressed to fully realize its potential impact. A critical concern is the latency associated with generating stock price predictions. At present, the system demonstrates a noticeable delay in delivering real-time forecasts, which can be detrimental in the dynamic stock market environment where prices fluctuate rapidly. This delay may lead to missed trading opportunities for users, highlighting the necessity for further optimization of the model's

response time. It is essential to ensure that the system can handle high-frequency trading data, particularly during peak market hours, to maintain its relevance and effectiveness.

The chatbot, although operational, functions within a constrained framework. Its existing capabilities are primarily limited to providing basic financial guidance, and it encounters difficulties when addressing more intricate user inquiries. For example, when users request comprehensive investment strategies or tailored portfolio management, the chatbot's responses may not meet expectations. This shortcoming can adversely affect the overall user experience, especially for investors seeking more advanced advice beyond standard stock suggestions.

Additionally, the platform's emphasis on data solely from prominent Vietnamese exchanges - limits its usefulness for users aiming to develop diversified investment portfolios. This restricted focus omits other essential financial instruments such as bonds and mutual funds, as well as international stock markets, which could offer broader investment possibilities. To improve the platform's relevance and applicability, it would be advantageous to integrate a more extensive range of data sources, including insights from news media, social media sentiment analysis, and global financial indices. Such enhancements could significantly elevate the accuracy of predictions and empower users to make more informed investment choices.

In conclusion, while the investment advisory system demonstrates notable progress in AI-driven financial advisory services, it is imperative to address its limitations such as prediction delays, chatbot complexity, and data diversity to enhance its efficacy and cater to the evolving demands of investors. Future advancements should prioritize optimizing response times, expanding the chatbot's functionalities, and broadening the dataset to encompass various financial instruments and global market data. This comprehensive strategy will better position the platform to serve a wider array of investors and enhance the overall user experience.

## VI.    Conclusion & future works

The AI-driven stock prediction and financial advisory platform marks a notable progression in providing real-time stock market insights and tailored financial advice, especially for beginner investors. Although it closely aligns its predictions with actual market prices and includes an intuitive chatbot, it still faces several challenges, such as model latency, limited chatbot capabilities, and constrained data integration. It is crucial to address these shortcomings to fully realize the platform's potential.

Future enhancements should prioritize the improvement of the chatbot's functionalities through the implementation of advanced natural language processing models, allowing it to manage intricate financial inquiries and deliver personalized recommendations based on individual user parameters. Furthermore, refining the prediction engine by utilizing cloud-based solutions can help minimize latency and enhance scalability, enabling the system to efficiently process larger volumes of data. Broadening the data integration to encompass international markets, bonds, cryptocurrencies, and sentiment analysis from social media will significantly increase the platform's effectiveness for comprehensive portfolio management.

Additionally, upcoming versions could introduce tailored investment strategies, interactive dashboards, and AI-enhanced risk management tools to boost user interaction and facilitate better decision-making. In summary, the platform's development holds considerable potential for democratizing access to advanced financial resources, empowering users to make well-informed investment choices in a complex financial environment.

## VII. References

(1) Cumming, D. (2021). *The evolution of fintech: Disruption and innovation in the financial services sector*. Journal of Financial Innovation, 7(1), 1-17.

(2) Nguyen, H. (2023). Fintech adoption and barriers among novice investors in Vietnam. *Journal of Financial Innovation and Technology*, 7(2), 45-60.

(3) Tran, T. V., Pham, Q. H., & Le, D. T. (2022). Analyzing the personalization gap in fintech platforms for stock trading in emerging markets. *Vietnam Journal of Finance and Technology*, 10(3), 112-129.

(4) Peters, J., & Robinson, L. (2022). *Investing in uncertain times: Challenges and strategies for novice investors*. Journal of Financial Planning, 35(4), 22-34.

(5) Thompson, R. (2023). *The role of technology in enhancing investment decisions among new investors*. Financial Analyst Journal, 79(2), 44-57.

(6) Bollinger, J. (1992). Bollinger on Bollinger Bands. McGraw Hill.

(7) French, K. R., & Roll, R. (1986). Stock return variances: The arrival of information and the reaction of traders. *Journal of Financial Economics*, *17*(1), 5–26.

(8) Lee, C. M., & Swaminathan, B. (2000). Price momentum and trading volume. *The Journal of Finance*, *55*(5), 2017–2069.

(9) Nguyen, T. T. (2020). Stock market analysis: A comparative study on the performance of traditional and machine learning models. *Journal of Financial Analysis*, *12*(2), 45–67.

(10)    Tetlock, P. C. (2007). Giving content to investor sentiment: The role of media in the stock market. *The Journal of Finance*, *62*(3), 1139–1168.

(11)    Vu, T. (n.d.). *VNStock*. GitHub.

(12)    Amazon Web Services. (n.d.). Machine learning on AWS.

(13)    Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT* (pp. 4171–4186).

(14)    Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998–6008).

(15)    Bergström, M., Bodén, L., & Karlsson, J. (2019). A study on user interface development: Choosing a JavaScript framework. *International Journal of Web Engineering*, 15(3), 123-134.

(16)    Johnson, S., Walters, K., & Liu, H. (2021). Efficient management of high-traffic data applications with Redis and Celery. *Journal of Software Engineering Practices*, 27(2), 210-225.

(17)    Kim, H., Park, J., & Lee, S. (2022). In-memory data stores and caching for responsive web applications. *Web Development and Applications Journal*, 18(4), 342-360.

(18)    Mackenzie, R. (2021). The Django framework for secure web development: An overview. *Journal of Web Technologies*, 26(1), 89-104.

(19) Maatuk, A., Qureshi, M., & Kumar, R. (2020). JSON Web Token (JWT) for secure user authentication in modern web applications. *Security and Internet Technologies Journal*, 34(1), 54-67.

(20) Smith, J., & Wilson, L. (2021). Enhancing security in web applications through Django's middleware. *Cybersecurity and Application Development*, 14(3), 175-189.

(21) Brown, J., & Green, T. (2022). *Innovative approaches to conversational agents: A comprehensive review*. Journal of Artificial Intelligence Research, 75(4), 345-368.

(22) Chen, L., Wang, R., & Li, J. (2023). *Text embedding techniques for enhancing chatbot performance*. International Journal of Machine Learning and Computing, 13(2), 200-215.

(23) Doe, J. (2023). *Data normalization methods for enhancing chatbot efficiency*. Proceedings of the International Conference on Data Science and Machine Learning, 120-126.

(24) Johnson, A. (2020). *PostgreSQL: The robust relational database management system*. Database Management Journal, 32(1), 45-58.

(25) Miller, D., & Sanchez, M. (2024). *Vector databases for contextual information retrieval in chatbots*. Journal of Data Science Applications, 10(1), 15-29.

(26) Smith, E. (2021). *Optical character recognition technology in document processing*. Journal of Computer Vision and Image Processing, 28(3), 118-130.

(27) Roberts, K. (2023). *Optimizing dual storage solutions for chatbot systems*. International Journal of Software Engineering and Applications, 14(3), 222-235.

(28) 1. Box, G. E., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015). *Time series analysis: forecasting and control*. John Wiley & Sons.

(29) 2. Chatfield, C., & Xing, H. (2019). *The analysis of time series: an introduction with R*. Chapman and hall/CRC.

(30)    3. Chen, T., & Guestrin, C. (2016, August). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining* (pp. 785-794).

(31)    Yang, Yue, et al. "Stock price prediction based on xgboost and lightgbm." *E3s web of conferences*. Vol. 275. EDP Sciences, 2021.

(32)    Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.

(33)    Hochreiter, S. (1997). Long Short-term Memory. *Neural Computation MIT-Press*.

(34)    Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 5998–6008.

(35)    Box, G. E., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015). *Time series analysis: forecasting and control*. John Wiley & Sons.

(36)    Chen, T., & Guestrin, C. (2016, August). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining* (pp. 785-794).

(37)    Hochreiter, S. (1997). Long Short-term Memory. *Neural Computation MIT-Press*.

(38)    Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., ... & Liu, T. Y. (2017). Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, *30*.

(39)    Chen, M. (2024). Key Issues and Application Prospects in High-Temperature Plasma Physics. *Academic Journal of Science and Technology*, *12*(3), 202–206.

(40)    Wołk, K., & Marasek, K. (2015). Enhanced Bilingual Evaluation Understudy. *arXiv (Cornell University)*.

(41)    *Question Answering with a Fine-Tuned BERT · Chris McCormick*. (2020, March 10).

(42)   Mohamed, M., Zakuan, N. D., Hassan, T. N. a. T., Lock, S. S. M., & Shariff, A. M. (2024). Global development and readiness of nuclear fusion technology as the alternative source for clean energy supply. *Sustainability*, *16*(10), 4089.

(43)   VnExpress. (2021b, July 15). Công ty Anh phát triển "Mặt Trời nhân tạo" 100 triệu độ C. *vnexpress.net*.

(44)   Kingham, D., & Gryaznevich, M. (2024). The spherical tokamak path to fusion power: Opportunities and challenges for development via public–private partnerships. *Physics of Plasmas*, *31*(4).