UNIVERSITY OF SCIENCE
**FACULTY OF INFORMATION TECHNOLOGY**


**Hồ Sỹ Kiên**


# ASSESSING THE IMPACT OF NEWS DATA ON STOCK TREND PREDICTION MODELS


BACHELOR OF SCIENCE IN COMPUTER SCIENCE

STANDARD PROGRAM


Ho Chi Minh City, 07/2025

UNIVERSITY OF SCIENCE

**FACULTY OF INFORMATION TECHNOLOGY**

**Hồ Sỹ Kiên - 21120091**

# ASSESSING THE IMPACT OF NEWS DATA ON STOCK TREND PREDICTION MODELS

BACHELOR OF SCIENCE IN COMPUTER SCIENCE

STANDARD PROGRAM

**THESIS ADVISORS**

MSc. Trần Văn Quý

Ho Chi Minh City, 07/2025

# Reassurances

I declare that the content of this thesis is my own research work under the guidance of MSc. Trần Văn Quý. All data and research results in this thesis are prensented honestly and do not overlap with other studies.

# Thesis Editing Explanation

# BẢN THUYẾT MINH CHỈNH SỬA BÁO CÁO ĐỀ TÀI
## KHÓA LUẬN TỐT NGHIỆP

Tên đề tài : Đánh giá tác động dữ liệu tin tức vào mô hình dự đoán xu hướng chứng khoán

Sinh viên thực hiện : ....Hồ Sỹ Kiên – 21120091 ........................................................

Hội đồng bảo vệ: ...Kỹ thuật phần mềm ..................................................................

Ngày bảo vệ: ...28/07/2025 ...................................................................................

(Chúng) tôi đã hoàn chỉnh khóa luận tốt nghiệp theo góp ý của Hội đồng và nhận xét của Giảng viên phản biện. Nội dung đã hiệu chỉnh như sau:

1. **Tên đề tài đã điều chỉnh:** *(nếu có)*
   - Tiếng Việt : ...............................................................................................
   .......................................................................................................................
   - Tiếng Anh : ...............................................................................................
   .......................................................................................................................

2. **Điều chỉnh nội dung báo cáo:** *(Ghi rõ nội dung, chỉ rõ vị trí chỉnh sửa: trang/chương...)*
   - Bổ sung kiến trúc tổng quan của mô hình PEN, mục tiêu của PEN, những hạn chế của mô hình và động lực cải tiến ở Chương 3, Phần 3.3 (từ trang 32 đến trang 39).
   - Cập nhật kết quả thí nghiệm của mô hình cải tiến ở Chương 4 (bảng 4.1, trang 50) và Chương 5 (bảng 5.1, trang 73).
   - Cập nhật và sửa lỗi chú thích cho hình kiến trúc DP-SRL bị tràn viền ở Chương 5 (hình 5.1, trang 81).
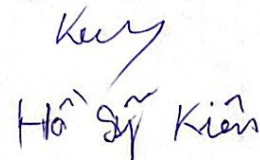
Tp. Hồ Chí Minh, ngày 14 tháng 8 ... năm .20.25

| **Xác nhận của Giảng viên hướng dẫn** | **(Nhóm) Sinh viên** |
|---|---|
| (Ký tên, ghi rõ họ tên) | (ký tên, ghi rõ họ tên) |

*Trần Văn Quý*

*Kiên*

*Hồ Sỹ Kiên*

**Xác nhận của Giảng viên**
**phản biện**
(Ký tên, ghi rõ họ tên)

*Nguyễn Thị Minh Quyền*

# Comment of Thesis's Advisor

4

# BẢN NHẬN XÉT KHÓA LUẬN TỐT NGHIỆP

## (HƯỚNG NGHIÊN CỨU)

*Tên đề tài : Đánh giá tác động dữ liệu tin tức vào mô hình dự đoán xu hướng chứng khoán*

*Sinh viên thực hiện :* **21120091 – Hồ Sỹ Kiên**

*Giảng viên hướng dẫn:* **ThS. Trần Văn Quý**

1. **Chủ đề và ý tưởng nghiên cứu:**

Nghiên cứu tác động của tin tức đến xu hướng chứng khoán, phân tích các yếu tố nhân quả của tin tức đối với biến động thị trường, và khảo sát các kỹ thuật kết hợp hiệu quả giữa dữ liệu giá và dữ liệu tin tức. Trên cơ sở đó, luận văn tiến hành thử nghiệm tích hợp các kết quả nghiên cứu để cải tiến mô hình PEN.

2. **Phương pháp nghiên cứu:**

Thu thập và tiền xử lý dữ liệu giá chứng khoán và tin tức tài chính; áp dụng các phương pháp phân tích nhân quả để xác định mối quan hệ giữa tin tức và xu hướng giá; thử nghiệm các kỹ thuật kết hợp dữ liệu giá và tin tức; tích hợp kết quả vào mô hình PEN và đánh giá hiệu quả thông qua các chỉ số dự báo và so sánh với các mô hình gốc.

3. **Đóng góp Khoa học và thực tiễn:**

Causal Text Selection Unit (Ca-TSU): Mô-đun lọc văn bản để tách các thông tin có ảnh hưởng nhân quả, tạo Vector of Causality (VoC) giúp tăng độ chính xác dự báo bằng cách tập trung vào tín hiệu dự đoán.

Dual-Pathway Shared Representation Learning (DP-SRL): Kiến trúc hai nhánh mô hình hóa mối quan hệ hai chiều giữa văn bản và giá, tách biệt văn bản mang tính dự đoán và văn bản mang tính phản ứng nhằm cải thiện độ chính xác và khả năng diễn giải.

Giảng viên hướng dẫn - 1

Volatility-Aware Fusion Mechanism: Cơ chế tích hợp động điều chỉnh trọng số giữa dữ liệu văn bản và giá dựa trên biến động thị trường, đảm bảo hiệu suất ổn định trong nhiều điều kiện khác nhau.

### 4. Quá trình thực hiện và quản lý dự án:
Sinh viên làm việc rất chăm chỉ nghiêm túc.

### 5. Báo cáo viết:
Báo cáo viết rõ ràng, đầy đủ, logic.

### 6. Trình bày trước hội đồng:
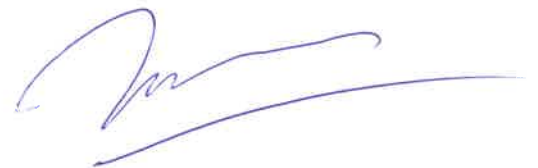Trình bày tốt.

### 7. Công bố khoa học/ ứng dụng thực tế:
Chưa có công bố khoa học.

*Đánh giá xếp loại:* Xuất sắc

TP.HCM, ngày 14 tháng 08 năm 2025
**Giảng viên hướng dẫn**
(*Ký và ghi rõ họ tên*)

Trần Văn Quý

Giảng viên hướng dẫn - 2

# Comment of Thesis's Reviewer

4

# BẢN NHẬN XÉT KHÓA LUẬN TỐT NGHIỆP

### (HƯỚNG NGHIÊN CỨU)

*Tên đề tài : Đánh giá tác động dữ liệu tin tức vào mô hình dự đoán xu hướng chứng khoán*
*Sinh viên thực hiện :        21120091 – Hồ Sỹ Kiên*
*Giảng viên phản biện: TS. Nguyễn Thị Minh Tuyền*

## 1. Chủ đề và ý tưởng nghiên cứu:

PEN (Prediction-Explanation Network) là một mô hình cho phép kết hợp luồng giá (price stream) và luồng văn bản (text tream) để đưa ra dự đoán về sự chuyển động của giá chứng khoán. Mục tiêu đề ra của đề tài này là nhằm cải tiến mô hình PEN để đạt được kết quả tốt hơn.

## 2. Phương pháp nghiên cứu:

Nhóm đã tiến hành tìm hiểu và đánh giá các công trình liên quan về dự đoán giá chứng khoán nói chung, và dự đoán giá có giải thích (Explainable stock prediction). Nhóm tiến hành đề xuất các module khác nhau như Ca-TSU (Causal Text Selection Unit) và DP-SRL (Dual-Path Shared Representation Learning) nhằm thay thế các module có sẵn trong mô hình. Nhóm bổ sung thêm cơ chế hợp nhất nhận thức biến động (Volatility-Aware Fusion Mechanism) nhằm tích hợp dữ liệu văn bản và giá dựa vào điều kiện thị trường hiện tại. Nhóm cũng tiến hành các thực nghiệm cần thiết để minh chứng cho việc cải thiện độ chính xác so với mô hình PEN ban đầu.

## 3. Đóng góp Khoa học và thực tiễn:

Dưới đây là một số đóng góp khoa học của nhóm:

- Nhóm đề xuất module Ca-TSU, thay cho module gốc là TSU. Module mới lọc các nội dung có liên quan đến việc dự đoán, giảm tác động của các văn bản không liên quan hoặc nhiễu, nhằm làm tăng khả năng diễn giải và tính hiệu quả của mô hình PEN.

- Module chính của mô hình PEN là SRL(Shared Representation Learning) có chức năng tìm hiểu văn bản nào có thể liên quan đến biến động giá cổ phiếu bằng cách mô hình hoá tương tác giữa dữ liệu văn bản và dữ liệu giá cổ phiếu mối tương quan của chúng. Nhóm đề xuất module thay thế DP-SRL nhằm phân biệt giữa thông điệp dự đoán (văn bản -> giá) và thông điệp phản ánh (giá -> văn bản), nắm bắt mối quan hệ hai chiều giữa giá và văn bản nhằm nâng cao khả năng dự đoán và khả năng diễn giải.

Giảng viên phản biện - 1

- Nhóm đề xuất module thứ 3 để hợp nhất dữ liệu văn bản và dữ liệu giá một cách linh hoạt dựa trên thông tin biến động thị trường, ưu tiên các đầu vào liên quan trong giai đoạn biến động hoặc ổn định.

Việc thực nghiệm, so sánh kết quả thực nghiệm với kết quả của mô hình gốc. Tuy nhiên, tác giả không đưa ra kết quả thực nghiệm của mô hình gốc trên cùng môi trường thực thi với Ca-TSU, chỉ sử dụng lại kết quả của bài báo gốc làm baseline. Điều này tương tự với DP-SRL. Module thứ 3 có giải thích kết quả thực nghiệm nhưng không cung cấp số liệu cụ thể.

**4. Báo cáo viết:**

Báo cáo viết có 118 trang, gồm 7 chương chính: Chương 1 giới thiệu ngữ cảnh thực hiện đề tài; Chương 2 trình bày kiến thức nền tảng liên quan đến các mô hình tạo sinh, VAE (Variational Autoencoder) ; Chương 2 giới thiệu về dự đoán giá chứng khoán, và dự đoán chứng khoán có giải thích; Chương 4, 5 tập trung vào hai đề xuất chính Ca-TSU và DP-SRL; Chương 6 đề cập đến cơ chế hợp nhất nhận thức biến động; Chương 7 kết luận và đề xuất hướng phát triển.

Nhìn chung, báo cáo có cấu trúc hợp lý. Tuy nhiên, việc trình bày còn thiếu tính logic dẫn đến việc khó nắm bắt được sự kết nối giữa các đóng góp trong báo cáo, sự liên quan giữa mô hình PEN với phần cải tiến. Khoá luận nên bắt đầu bằng việc trình bày kiến trúc của mô hình PEN để người đọc có kiến thức ban đầu, sau đó chỉ ra các phần được cải tiến. Nhóm sinh viên cũng nên bổ sung thêm các hình vẽ để minh hoạ, thay vì giải thích hoàn toàn bằng văn bản.

**5. Trình bày trước hội đồng:**

Tương đối tốt, sinh viên cần nắm rõ bài toán đặt ra, ở mỗi đề xuất nên hiểu rõ mình đang giải quyết vấn đề gì, nhằm mục đích gì, và vì sao mình lại thực hiện như vậy.

**6. Công bố khoa học/ ứng dụng thực tế:**

Không có công bố khoa học.

*Đánh giá xếp loại: Giỏi-Xuất sắc*

TP.HCM, ngày 28 tháng 07 năm 2025

**Giảng viên phản biện**

(*Ký và ghi rõ họ tên*)

**Nguyễn Thị Minh Tuyền**

Giảng viên phản biện - 2

# Acknowledgement

First and foremost, I would like to express my deep appreciation to my thesis advisor, MSc.Trần Văn Quý, for his dedicated supervision, insightful guidance, and constructive feedback throughout the duration of this thesis. His extensive knowledge, expertise, and dedication have shaped the direction and ensured the quality of this thesis.

I am also grateful to the Faculty of Information Technology, University of Science – VNU-HCM, for providing a rigorous academic environment and the necessary resources to complete this thesis sucessfully. The academic atmosphere has genuinely enriched my academic experience.

I would like to extend my sincere thanks to my friends' support and encouragement. Their belief in our abilities, willingness to provide insightful feedback, and constant inspiration and motivation have been invaluable throughout this journey. Their presence has made a significant difference, and I am grateful for their unwavering friendship and support.

To my family, I am profoundly thankful for your unwavering support, patience, and belief in my academic pursuits. Their support and understanding have been a constant source of strength, and I am forever grateful for their presence in my life.

Finally, I acknowledge with gratitude all individuals who have, in various ways, supported and contributed to the realization of this research. Their contributions, whether big or small, have played a significant role in shaping the outcome of this thesis. I am truly grateful for their involvement and support.

# Thesis Outline

GRADUATION THESIS OUTLINE

# ASSESSING THE IMPACT OF NEWS DATA ON STOCK TREND PREDICTION MODELS

*(Đánh giá tác động dữ liệu tin tức vào mô hình dự đoán xu hướng chứng khoán)*

# 1 GENERAL INFORMATION

**Instructor:**

– MSc. Trần Văn Quý (Faculty of Information Technology)

**Student:**

1. Hồ Sỹ Kiên (Student ID: 21120091)

**Project type:** Research

**Project timeline:** From *01/2025* to *07/2025*

# 2 THESIS CONTENT

## 2.1 Introduction

Stock price prediction remains a cornerstone of financial analysis, yet it is an inherently complex task due to the multifaceted influences on market dynamics. Traditional models predominantly rely on historical price data and technical indicators to forecast future stock movements. However, these approaches often fall short in capturing the full spectrum of factors driving stock prices, such as economic events, corporate developments, and shifts in market sentiment. News data—encompassing financial articles, social media posts, and other textual sources—offers a rich, real-time window into these external influences. Assessing the impact of news data involves evaluating how effectively these sources can be integrated into prediction models to enhance accuracy, precision, and interpretability, thereby aiding investors, traders, and financial analysts in decision-making.

The integration of news data is particularly relevant in the era of big data and advanced machine learning, where natural language processing (NLP) techniques enable the extraction of meaningful information from text. Recent research ([1, 2, 3]) suggests that news data can capture market reactions to events before they are fully reflected in price data, potentially improving prediction models. Similarly, hybrid architectural approaches ([4, 5, 6]) have shown promising results. These findings highlight the potential of combining different modeling paradigms to capture complex market dynamics. However, challenges include processing large volumes of text, handling noise, and ensuring the model remains interpretable, especially in finance where transparency is crucial for regulatory compliance and trust.

In this thesis, I aim to study and develop a hybrid model that synergistically combines traditional price-based analysis with insights derived from news data. The core idea is to leverage natural language processing (NLP) techniques to

extract meaningful information—such as sentiment and thematic content—from news sources and integrate it with conventional stock price data. This approach aims to create a more holistic prediction framework that accounts for both historical patterns and real-time external influences, ultimately improving the model's predictive power and contextual relevance.

The practical significance of this research extends across multiple domains. For individual investors, more accurate stock price predictions can inform better investment decisions, while financial institutions stand to benefit from improved risk management and portfolio optimization. Beyond these direct applications, understanding the interplay between news data and stock prices offers broader insights into market behavior and investor sentiment, providing value to economists and policymakers. In an era defined by information overload, the ability to efficiently filter and utilize relevant news data also addresses the growing demand for transparency and explainability in automated financial systems, aligning with regulatory expectations and fostering greater trust in predictive technologies.

## 2.2   Objectives

I am dedicated to researching and developing a prediction-explanation system with framework built upon hybrid models, refining the integration of news data to capture market sentiments and events, thereby improving predictive power while providing clear insights into the underlying causes of those predictions. In the context of financial markets, where stock price prediction traditionally relies on historical price data and technical indicators, this approach can help addressing a critical limitation: the lack of transparency and interpretability in many advanced machine learning models.

The proposed approach has serveral potential benefits. It enhances prediction accuracy by leveraging external information beyond historical data and improves explainability by identifying specific news items or sentiments that impact stock

movements. The possible impacts are substantial for individual investors and financial institutions, more accurate and transparent models can lead to better investment decisions, improved risk management, and increased confidence in automated systems, thus advancing both practical applications and theoretical research directions in finance and beyond.

## 2.3 Scope

The scope of the thesis will involve conducting a series of studies to explore and evaluate the capability of explainable AI, particularly hybrid stock price prediction model that integrates news data from a social network with traditional price-based analysis, aiming to enhance both predictive accuracy and explainability.

## 2.4 Expected Approach

**Hybrid Stock Price Prediction Model:** Hybrid models combine multiple data sources and methodologies to enhance predictive accuracy and provide explainable insights into stock price movements. This thesis focuses on integrating news data from social networks with traditional price-based analysis, leveraging advanced machine learning techniques to capture market sentiments and events. Hybrid architectures have emerged as cutting-edge solutions for financial forecasting, excelling in processing heterogeneous data sources, identifying complex relationships, and generating interpretable predictions. This advancement has led to applications such as sentiment-driven trading strategies, event impact analysis, and explainable AI systems for financial decision-making. Notably, hybrid models have been employed in multi-modal prediction tasks. Several approaches, including Deep Fusion Model ([7, 8, 6]), LSTM-GNN ([4, 5]), Transformer-based sentiment analysis([9, 10]), and explainable AI frameworks ([1, 11, 12, 13]), have demonstrated significant improvements in predictive accuracy and transparency.

News Data Integration: The integration of news data aims to capture senti-

ment and event-driven market dynamics that traditional price-based models may overlook. News data from financial articles, social media platforms, and company reports is processed using Natural Language Processing (NLP) techniques to extract sentiment scores, topic relevance, and event impact indicators. This process involves advanced text representation methods such as embeddings (e.g., Word2Vec, BERT), attention mechanisms, and transformer architectures. Text serves as crucial input for hybrid models, enabling the identification of contextual relationships between market events and price movements. By embedding sentiment-driven insights into predictive frameworks, this approach expands the scope of traditional stock prediction models to incorporate real-time market reactions.

**Explainability Framework:** The explainability component aims to provide transparent insights into the predictions generated by the hybrid model. This involves developing a Summarize-Explain-Predict (SEP) ([12]) framework that generates human-readable explanations for stock price movements based on contributing factors such as news sentiment and social network diffusion patterns. Techniques such as attention visualization, counterfactual analysis, and feature attribution methods such as SHAP, LIME ([14]), VoS ([1]) are employed to enhance interpretability. Explainability serves as a crucial aspect of the thesis, ensuring that predictions are actionable and aligned with investor expectations.

**My approach:** My approach focuses on studying the capability of explainable AI, particularly hybrid stock price prediction model with news data integration and an effective explainability framework, thus enhancing the ability to contribute to investment decisions. By leveraging advanced NLP techniques, the model captures market sentiments and event-driven fluctuations from financial news and social media. These insights are then combined with historical price data using deep-learning architectures to uncover complex relationships. Additionally, I incorporate the explainability framework to provide interpretable predictions, with

an emphasis on both transparency and explainability.

## 2.5 Expected Results

**Experimental results:** Through the experiments conducted throughout the thesis, we can draw the advantages and disadvantages of the proposed model building methods (news integrability, explainability, etc.). From these, possible directions for improvements ban be suggested in the future.

## 2.6 Thesis Plan

| From | To | Task |
|---|---|---|
| 01/01/2025 | 31/01/2025 | Survey Hybrid Stock Price Prediction Model |
| 01/01/2025 | 31/01/2025 | Survey News Data Integration |
| 01/01/2025 | 31/01/2025 | Survey Explainability Framework |
| 01/02/2025 | 01/03/2025 | Study Hybrid Stock Price Prediction Model |
| 01/02/2025 | 01/03/2025 | Study News Data Integration |
| 01/02/2025 | 01/03/2025 | Study Explainability Framework |
| 02/03/2025 | 31/03/2025 | Employ Hybrid Stock Price Prediction Model |
| 02/03/2025 | 31/03/2025 | Employ News Data Integration |
| 02/03/2025 | 31/03/2025 | Employ Explainability Framework |
| 01/04/2025 | 31/05/2025 | Improve Hybrid Stock Price Prediction Model |
| 01/04/2025 | 31/05/2025 | Improve News Data Integration |
| 01/04/2025 | 31/05/2025 | Improve Explainability Framework |
| 01/06/2025 | 31/06/2025 | Develop prototype |
| 01/07/2025 | 31/07/2025 | Conduct experiments |

# References

[1] S. Li, W. Liao, Y. Chen, and R. Yan, "Pen: Prediction-explanation network to forecast stock price movement with better explainability," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 4, pp. 5187–5194, 2023.

[2] D. Hirshleifer, L. Peng, and Q. Wang, "News diffusion in social networks and stock market reactions," *The Review of Financial Studies*, vol. 38, pp. 883–937, March 2025.

[3] B. A. Neri-Mares, V. A. Rodriguez-Ríos, R. R. Gallegos-Villela, and E. J. Suarez-Dominguez, "Influence of social media on the stock market: Part 1. a brief analysis," *Universal Journal of Business and Management*, vol. 4, no. 1, pp. 1–14, 2024.

[4] M. S. Sonani, A. Badii, and A. Moin, "Stock price prediction using a hybrid lstm-gnn model: Integrating time-series and graph-based analysis," *arXiv preprint*, vol. 2502.15813, February 2025.

[5] K. Yadav, M. Yadav, and S. Saini, "Stock values predictions using deep learning based hybrid models," *CAAI Transactions on Intelligence Technology*, vol. 7, pp. 107–116, June 2021.

[6] L. Liu, "Presenting an optimized hybrid model for stock price prediction," *International Journal of Advanced Computer Science and Applications (IJACSA)*, vol. 15, no. 1, 2024.

[7] P. Chen, Z. Boukouvalas, and R. Corizzo, "A deep fusion model for stock market prediction with news headlines and time series data," *Neural Computing and Applications*, vol. 36, pp. 21229–21271, August 2024.

[8] R. Corizzo and J. Rosen, "Stock market prediction with time series data and news headlines: a stacking ensemble approach," *Journal of Intelligent Information Systems*, vol. 62, pp. 27–56, 2024.

[9] W. Gu, Y. Zhong, S. Li, C. Wei, L. Dong, Z. Wang, and C. Yan, "Predicting stock prices with finbert-lstm: Integrating news sentiment analysis," *arXiv preprint arXiv:2407.16150*, 2024.

[10] M. Patel, K. Jariwala, and C. Chattopadhyay, "A hybrid relational approach toward stock price prediction and profitability," *IEEE Transactions on Artificial Intelligence*, vol. 5, pp. 5844–5854, Nov. 2024.

[11] K. Du, R. Mao, F. Xing, and E. Cambria, "Explainable stock price movement prediction using contrastive learning," in *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management (CIKM '24)*, (New York, NY, USA), pp. 529–537, ACM, 2024.

[12] K. J. L. Koa, Y. Ma, R. Ng, and T.-S. Chua, "Learning to generate explainable stock predictions using self-reflective large language models," *arXiv preprint arXiv:2402.03659*, 2024.

[13] S. M. Carta, S. Consoli, L. Piras, A. S. Podda, and D. R. Recupero, "Explainable machine learning exploiting news and domain-specific lexicon for stock market forecasting," *IEEE Access*, vol. 9, pp. 30193–30205, 2021.

[14] D. Muhammad, I. Ahmed, K. Naveed, and M. Bendechache, "An explainable deep learning approach for stock market trend prediction," *Heliyon*, vol. 10, no. 1, p. e12345, 2024.

| | |
|---|---|
| **XÁC NHẬN**<br>**CỦA NGƯỜI HƯỚNG DẪN**<br>*(Ký và ghi rõ họ tên)* | *TP. Hồ Chí Minh, ngày… tháng… năm…*<br>**NHÓM SINH VIÊN THỰC HIỆN**<br>*(Ký và ghi rõ họ tên)* |

# Contents

# List of Figures

# List of tables

# Abstract

As financial institutions navigate increasingly complex markets and regulatory environments, the demand for interpretable and accountable stock price prediction models has intensified. Regulatory frameworks such as the EU's MiFID II and emerging AI governance guidelines now require not only high predictive accuracy but also transparent, explainable decision processes—particularly when models influence investment or trading actions. In response, researchers have turned to machine learning models that incorporate textual information from financial news and social media, leveraging the assumption that public sentiment and breaking news significantly influence market dynamics.

However, a key limitation persists in these approaches: they often fail to distinguish between texts that precede and potentially cause market movements and those that merely describe or react to events after they have occurred. This conflation can introduce bias, reduce predictive power, and obscure causal inference, undermining both performance and interpretability.

To address this challenge, I propose the Causal and Dual-Pathway Prediction-Explanation Network (CDP-PEN), a novel model architecture designed to enhance both the predictive performance and interpretability of text-based stock forecasting by explicitly separating predictive signals from reactive noise. CDP-PEN introduces two key components. First, the Causal Text Selection Unit (Ca-TSU) identifies and filters text embeddings to isolate those with predictive value. This process yields a Vector of

Causality (VoC), where each scalar quantifies the causal importance of individual textual inputs. Second, the Dual-Pathway Shared Representation Learning (DP-SRL) module explicitly models the bidirectional relationship between text and stock prices. This module incorporates pathway-specific attention mechanisms to capture predictive signals, gated communication to enable inter-pathway information flow, and volatility-aware fusion gates to dynamically weight the relevance of each pathway.

Empirical evaluations on real-world financial datasets demonstrate that CDP-PEN surpasses existing state-of-the-art models in predictive accuracy. Moreover, its structured dual-pathway design enhances interpretability by disentangling causally predictive content from reactive narratives. This improved explainability supports financial analysts in understanding market behavior more comprehensively and enables the assessment of informational value under varying market conditions, thereby fostering more informed decision-making.

# Chapter 1

# Introduction

*This chapter presents our motivations and objectives for exploring how generative AI can enhance stock price movement forecasting and fostering informed decision-making. We summarize our main contributions and describe the overall structure and main content chapters of our thesis.*

## 1.1 Overview

Financial markets operate as intricate systems where information, particularly news, plays a pivotal role in influencing price movements. The impact of news on these markets is well-documented, with unexpected news events often driving significant shifts in asset prices [1]. This phenomenon, known as the "announcement effect," highlights how public announcements, especially from governmental or monetary authorities, can directly affect market volatility and security prices [2]. Research utilizing text-mining techniques has shown that the sentiment expressed in news articles significantly influences asset prices, with foreign news often exerting a more pronounced effect on local markets than domestic news [3], [4], [5]. For instance, a study analyzing over 4 million Reuters articles from 1991 to 2015 found that sudden changes in news sentiment correlate with international asset price movements, with a global news-based sentiment index outper-

forming traditional indicators like the CBOE Volatility Index in predicting market trends [3]. Furthermore, network analysis reveals that sentiment events in news can propagate through financial networks, affecting not only individual companies but also groups of related firms, thereby amplifying market movements [6].

The advent of machine learning (ML) has transformed the landscape of financial markets by addressing complex challenges through data-driven solutions. ML models are now integral to various financial applications, including stock market prediction, risk management, fraud detection, and algorithmic trading [7]. These models excel in analyzing vast datasets to identify patterns and make predictions, thereby enhancing decision-making processes. For example, deep learning techniques, such as long short-term memory (LSTM) networks, have been employed for stock price prediction, often outperforming traditional methods [8]–[10]. The integration of ML into financial markets underscores its potential to enhance efficiency and accuracy in decision-making, though it also necessitates robust frameworks for oversight and risk management.

In recent years, the application of generative artificial intelligence (AI) has emerged as a critical tool for increasing the explainability of complex financial models. Explainability, or the ability to understand how AI models arrive at their decisions, is essential for building trust and ensuring transparency in financial services. Generative AI enhances this by providing insights into the decision-making process, allowing stakeholders to comprehend the inputs and data used to generate recommendations. This is particularly important in finance, where regulatory compliance and ethical considerations demand that AI-driven decisions be interpretable and unbiased [11]. For example, generative AI has been used to assess the impact of regulatory changes, such as new capital rules, thereby aiding in the interpretation of complex regulatory frameworks [12]. Additionally, generative AI's ability to process unstructured data, such as historical service interactions or news articles, can democratize access to insights,

further supporting explainability in financial operations. Specific studies have demonstrated the efficacy of generative models in this domain. For instance, Generative Adversarial Networks (GANs) have been applied to generate synthetic financial time series data, capturing stylized facts such as random walks and time-varying volatility, which can be used to enhance the interpretability of risk management models by simulating realistic market scenarios [13]–[15]. Similarly, Variational Autoencoders (VAEs) have been utilized to create synthetic financial datasets for training predictive models, improving risk assessment by providing transparent data representations [16], [17]. Large Language Models (LLMs), such as those based on transformer architectures, have been employed to analyze financial news and reports, extracting interpretable insights for investment decision-making and market trend predictions [18] However, the deployment of generative AI must be accompanied by robust governance mechanisms to mitigate risks such as model "hallucinations" and ensure interpretability [12].

Building on these observations, my work focuses on developing techniques to advance the abilities of explainable generative models in stock price movement prediction.

## 1.2   Motivations

The increasing complexity and dynamism of financial markets underscore the need for advanced analytical tools capable of addressing the limitations of traditional models and harnessing emerging technologies to improve decision-making processes. Traditional financial models, such as autoregressive integrated moving average (ARIMA) and other time-series forecasting techniques, primarily rely on historical price data to predict future market trends [19]. While these models have been effective in certain contexts, their dependence on structured numerical data limits their ability to incorporate diverse information sources, such as news sentiment or macroeconomic events, which are critical drivers of market behavior.

3

Moreover, traditional models often lack explanatory power, which is known as "black-box" problem, providing little insight into the underlying factors driving their predictions. This opacity poses significant challenges in financial applications, where stakeholders, including regulators and investors, demand transparency to ensure trust and compliance with ethical and regulatory standards.

Recent advances in machine learning (ML) and, more notably, generative artificial intelligence (AI) has offered promising avenues to overcome these limitations. However, current generative models, such as Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs), and Large Language Models (LLMs), face their own set of challenges that hinder their full potential in financial applications. One significant issue is their susceptibility to noise in input data, particularly when processing unstructured sources like news articles or social media, which can introduce biases or inaccuracies in model outputs. For instance, GANs used to generate synthetic financial time series may struggle to capture the nuanced stochastic properties of market data, leading to unreliable simulations. Similarly, LLMs, while adept at extracting insights from textual data, can produce "hallucinations" outputs that are plausible but factually incorrect—posing risks in high-stakes financial decision-making. Additionally, current generative models typically treat all text data as having the same relationship with price movements, failing to distinguish between texts that predict future movements and those that merely respond to past changes. Finally, these models often lack the ability to adapt to varying market contexts, where the relative importance of price and text information fluctuates.

These challenges highlight the need for innovative approaches that combine the strengths of generative AI with mechanisms to enhance robustness, interpretability, and adaptability. Developing models that can effectively integrate diverse data sources, mitigate noise, and provide clear explanations for their outputs is critical for advancing their application in

4

financial markets. Such advancements would not only improve predictive performance but also align with regulatory requirements for transparency and accountability, thereby fostering greater trust among stakeholders. Research into hybrid models that leverage the generative capabilities of GANs, VAEs, or LLMs alongside traditional ML techniques could address these gaps, enabling more accurate and interpretable financial forecasting. This motivation drives ongoing efforts to refine AI methodologies, ensuring they meet the rigorous demands of financial applications while navigating the complexities of market dynamics.

## 1.3   Objectives and Main Contributions

With these motivations in mind, this thesis aims to develop an AI-based solution named **C**ausal and **D**ual-**P**athway **P**rediction-**E**xplanation **N**etwork (**CDP-PEN**) across three core dimensions of financial forecasting: predictive modeling, explainability, and robustness under market volatility. Together, these objectives target key limitations in current AI applications within finance, offering practical and theoretically grounded improvements.

First, while existing models often process textual and numerical data separately, there remains a gap in effectively combining these modalities to enhance forecasting accuracy. This work addresses that gap by proposing a unified model that integrates financial news and price history, allowing for a more comprehensive understanding of market signals and improving stock movement prediction performance.

Second, explainability is a growing concern in AI-driven financial systems, especially under regulatory and operational scrutiny. To meet this need, the proposed approach incorporates structured mechanisms to isolate causally relevant textual information, enabling clearer attribution of predictive outcomes. This contributes not only to model transparency but also to its usability in decision-critical environments.

Lastly, financial markets are inherently volatile and subject to abrupt regime changes, often triggered by external events. Most existing models exhibit performance degradation under such conditions. In response, this thesis introduces architectural features designed to maintain stable predictive capabilities during periods of heightened uncertainty, ensuring resilience and adaptability across a range of real-world scenarios.

I make the following main contributions:

- To improve the quality and relevance of textual signals in stock prediction, I propose the **Causal Text Selection Unit** (**Ca-TSU**), a dedicated module that identifies and filters textual inputs with causal influence on market movements. By isolating predictive messages from irrelevant noise, **Ca-TSU** enhances the model's focus on meaningful content and supports more accurate forecasting.

- To capture the complex interaction between market narratives and price behavior, I introduce **Dual-Path Shared Representation Learning** (**DP-SRL**). This dual-pathway architecture explicitly separates predictive texts (which precede price changes) from reflective texts (which respond to them), enabling the model to learn bidirectional dependencies between text and price. This separation improves not only prediction performance but also interpretability by highlighting how different types of text contribute to market dynamics.

- To ensure stability across varying financial environments, I design a **Volatility-Aware Fusion Mechanism** that dynamically adjusts how textual and numerical information is integrated. By accounting for market volatility , this mechanism enables the model to maintain robustness and adaptability under diverse and unpredictable conditions.

Extensive experiments on real-world financial datasets demonstrate that CDP-PEN model achieves superior accuracy compared to state-of-

6

the-art baselines while providing significantly enhanced explainability. By explicitly modeling the causal and reactive roles of textual information and dynamically adapting to market volatility, the model offers a transparent view of how different inputs influence stock price movements. Furthermore, the ability to trace information flow across pathways equips investors and analysts with deeper insights into market behavior and supports a more informed assessment of the relative impact of various information sources under shifting market conditions.

## 1.4 Thesis Organization

My thesis is structured as follows:

**Chapter 2** provides the theoretical foundations for stock price prediction, with a particular emphasis on Variational Autoencoders (VAEs).

**Chapter 3** reviews the current state-of-the-art in stock price prediction, with a focus on explainable AI (xAI). Specifically, I examine the latest advancements in VAEs models, considering their applications and contributions to the problem I am addressing.

**Chapter 4** introduces the **C**ausal **T**ext **S**election **U**nit (**Ca-TSU**), the first major contribution of this thesis. This module filters textual inputs, such as news and social media, to isolate messages with causal influence on stock prices. By employing a specialized attention mechanism, the Ca-TSU prioritizes predictive signals, reducing the impact of irrelevant or noisy text. The chapter details the implementation within the proposal framework, including the configuration of attention weights and regularization strategies. Experimental results demonstrate the effectiveness of the Ca-TSU in enhancing prediction accuracy by focusing on high-impact textual data, thereby addressing the challenge of noise in unstructured inputs.

**Chapter 5** presents the **D**ual-**P**ath **S**hared **R**epresentation **L**earning (**DP-SRL**) module, the second key contribution. This component distinguishes between predictive (text $\rightarrow$ price) and reflective (price $\rightarrow$ text)

messages through a dual-pathway architecture. By maintaining separate causal and responsive pathways, the DP-SRL captures the bidirectional relationship between text and price data, enhancing both predictive power and interpretability. The chapter describes the cross-pathway communication mechanism, implemented with adaptive gating and projection matrices, and evaluates its performance in modeling complex text-price interactions. Comprehensive experiments validate the module's ability to provide transparent insights into the factors driving market predictions.

**Chapter 6** focuses on the **Volatility-Aware Fusion Mechanism**, the third major contribution. This mechanism dynamically integrates text and price data based on current market conditions, prioritizing relevant inputs during volatile or stable periods. Additionally, the volatility-aware fusion employs context-dependent gates to weigh the contributions of price, causal text, and responsive text. The chapter discusses the technical details, including the normalization of gating mechanisms and their integration with the dual-pathway architecture. Experimental evaluations highlight the mechanism's robustness in adapting to rapid market fluctuations, ensuring consistent performance across diverse financial scenarios.

**Chapter 7** summarizes the thesis, synthesizing the key findings and contributions of my proposal model. It reflects on the model's advancements in prediction accuracy, interpretability, and adaptability, and discusses their implications for financial analytics. The chapter also outlines potential limitations and proposes directions for future research. This concluding chapter underscores the transformative potential of the proposed framework in addressing the evolving challenges of financial market prediction.

# Chapter 2

# Background

*This chapter presents the foundational background underpinning the proposed work. It begins with an overview of generative models, outlining the core principles and characteristics of the major model families. Particular emphasis is placed on Variational Autoencoders (VAEs), a prominent class of generative models that have garnered increasing attention in financial forecasting tasks for their capacity to generate high-quality, realistic data representations, making them particularly well-suited for modeling complex, stochastic processes such as stock price movements.*

## 2.1   Generative Models

A generative model is a probabilistic framework designed to approximate the underlying distribution from which observed data are drawn. Given a dataset $\mathcal{D}$ consisting of a finite number of samples assumed to originate from an unknown underlying distribution $p_{\text{data}}$, the primary objective of a generative model is to learn and replicate the statistical characteristics and structural patterns inherent in $p_{\text{data}}$. This enables the model to generate new data instances that are consistent with the original distribution.

### 2.1.1 Learning a generative model

Our focus is on parametric approximations of the data distribution, wherein the entire information contained in the dataset $\mathcal{D}$ is captured by a finite set of parameters. The objective in learning a generative model is to identify the parameter vector $\theta$ within a model family $\mathcal{M}$ such that the model distribution $p_\theta$ closely approximates the true data distribution $p_{\text{data}}$. This objective can be formally expressed as the following optimization problem:

$$\min_{\theta \in \mathcal{M}} d(p_{\text{data}}, p_\theta), \tag{2.1}$$

where $d(\cdot)$ denotes a suitable divergence or distance metric between probability distributions (e.g., Kullback–Leibler divergence, Jensen–Shannon divergence, or Wasserstein distance). An illustration of this approximation process is provided in Figure 2.1.



Figure 2.1: Finding a set of parameters $\theta$ within a model family $\mathcal{M}$ that minimizing the distance function $d(\cdot)$. Image source: Class note from Stanford CS236.

### 2.1.2 Applications of generative models

Generative models serve as a foundation for a range of downstream inference tasks, including sampling, density estimation, and unsupervised representation learning. However, it is important to recognize that not all

generative model classes support efficient or accurate inference across these tasks simultaneously. This inherent variability in task-specific performance has led to the development of diverse generative modeling approaches, each optimized for particular inference objectives.

**Sampling** is a central capability of generative models, aimed at generating novel and previously unseen data points $x_{\text{new}} \sim p_\theta(x)$ that are statistically consistent with the training distribution. This enables generative models to produce diverse data samples that reflect the patterns, structures, and variability present in the training dataset, supporting applications such as image synthesis and text generation

**Density estimation** involves computing the likelihood of new data points under the learned model distribution. For example, if a generative model $p_\theta(x)$ is trained on historical stock price movements conditioned on financial news, it should assign high probability to price patterns consistent with market norms and low probability to anomalous or rare movements. Accurate density estimation supports tasks such as outlier detection, risk assessment, and identifying events that deviate from typical market behavior.

**Unsupervised representation learning** is concerned with discovering meaningful and compact representations of high-dimensional data without reliance on labeled examples. The primary objective is to uncover the latent structure of the data and to extract salient features that encapsulate the most informative and relevant characteristics. These learned representations facilitate a broad range of downstream applications, including clustering, low-dimensional visualization, anomaly detection, and the enhancement of performance in supervised learning tasks by serving as informative feature embeddings.

### 2.1.3 Types of generative models

By varying the components of the optimization framework in Equation 2.1, the representation of the model family $\mathcal{M}$, the choice of divergence function $d(\cdot)$, and the optimization strategy employed to minimize $d(\cdot)$, a wide array of generative model families can be constructed.

In particular, **deep generative models** leverage deep neural networks to parameterize the model family $\mathcal{M}$, enabling them to capture highly complex data distributions. Deep generative models encompass a variety of architectures and learning approaches, each with its own strengths and characteristics. Notable classes of deep generative models include Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs), Normalizing Flows, and Diffusion Models, each of which embodies distinct architectural principles and learning paradigms. Figure 2.2 provides an overview of these architectures. Specifically:

- **GANs** utilize adversarial training, wherein a generator network produces synthetic data, while a discriminator network attempts to distinguish between real and synthetic data. This minimax optimization promotes the generation of realistic samples, useful in data augmentation and scenario stress testing.

- **VAEs** adopt an encoder-decoder framework to approximate the data likelihood by maximizing the Evidence Lower Bound (ELBO).

- **Normalizing Flows** learn an invertible transformation between a simple base distribution (e.g., Gaussian) and a complex target distribution.

- **Diffusion Models** define a forward Markov process that gradually adds noise to the data and a learned reverse process that reconstructs clean samples from noise.

Figure 2.2: Overview of leading deep generative model types. Image source: [20]

A successful generative model must balance three fundamental objectives: *sample quality*, *sample diversity*, and *sampling speed*. These dimensions are commonly referred to as the **generative learning trilemma**. Figure 2.3 illustrates how various generative model families perform relative to these criteria. Notably, no single model class currently achieves optimal performance across all three dimensions, highlighting ongoing challenges and research opportunities in the development of generative models tailored to financial prediction tasks.



Figure 2.3: Trade-offs across quality, diversity, and efficiency: the generative learning trilemma. Image source: [21]

## 2.2 Variational Autoencoders Models

A Variational Autoencoder (VAE) [22]–[26] is a generative model that extends the traditional autoencoder framework by introducing a probabilistic approach to learn a latent representation of data, enabling the generation of new samples. Its core components include an encoder that maps input data to a latent space with a distribution (typically Gaussian), parameterized by mean and variance, and a decoder that reconstructs data from samples drawn from this latent distribution, optimized via a loss function combining reconstruction error and a regularization term (KL divergence) to enforce a structured latent space. The illustration for the idea of diffusion models is depicted in Figure 2.4



Figure 2.4: The encoder and decoder process in VAEs models. Image source: [27]

### 2.2.1 Encoder Process

Starting with an observed data point $x$ (e.g. a sequence of stock prices) from the real data distribution, the VAE's *encoder* (also called the recognition or inference model) maps $x$ to a latent representation $z$ in a probabilistic manner [22], [28]. In a standard VAE, this mapping is implemented as a neural network that outputs the parameters of an approximate posterior distribution $q_\phi(z \mid x)$. A common choice is to assume a factorized

15

Gaussian form for the latent distribution, according to Aymne and Bernstein [23], [24]. In other words, instead of encoding $x$ to a single point in latent space, the encoder produces a mean vector $\boldsymbol{\mu}_\phi(x)$ and variance vector $\boldsymbol{\sigma}_\phi^2(x)$, defining

$$q_\phi(z \mid x) = \mathcal{N}\Big(z; \boldsymbol{\mu}_\phi(x), \operatorname{diag}\big(\boldsymbol{\sigma}_\phi^2(x)\big)\Big). \tag{2.2}$$

This probabilistic encoding allows the model to capture uncertainty and variability in the data [26]. Because sampling from $q_\phi(z \mid x)$ is part of the forward pass during training, Kingma and Max [22] employ the *reparameterization trick* . Concretely, letting $\varepsilon \sim \mathcal{N}(0, I)$, it could be set

$$z = \boldsymbol{\mu}_\phi(x) + \boldsymbol{\sigma}_\phi(x) \odot \varepsilon, \tag{2.3}$$

so that $z$ becomes a deterministic function of $x$ and the auxiliary noise $\varepsilon$, enabling gradients to propagate through the sampling operation. The encoder thus provides an *amortized inference* mechanism: a single set of parameters $\boldsymbol{\phi}$ is used to infer latent distributions for *all* data points, rather than optimizing separate variational parameters for each sample [29]. By producing a distribution in latent space (rather than a point estimate), the encoder ensures that the representation of data remains smooth and meaningful, which is crucial for generative modeling and downstream forecasting tasks. Although the above assumes continuous $z$, the framework also handles *discrete* latents by outputting class probabilities and resorting to special gradient estimators (see Section 2.2.3) [28].

## 2.2.2 Decoder Process

The *decoder* (generative model) strives to reconstruct or generate data $x$ from a latent code $z$. The generative process begins by drawing a latent vector from the prior distribution $p(z) = \mathcal{N}(0, I)$, which both regularizes the latent space and simplifies inference [23], [27]. Given a sample $\mathbf{z}$, the

16

decoder produces the parameters of the likelihood distribution for $\mathbf{x}$, which could be defined (according to Class note from Stanford CS236):

$$p_{\boldsymbol{\theta}}(x, z) \;=\; p(z)\, p_{\boldsymbol{\theta}}(x \mid z) \tag{2.4}$$

For continuous data such as stock prices, a common choice is to model $p_{\theta}(\mathbf{x}|\mathbf{z})$ as a Gaussian whose mean is given by a neural network $f_{\theta}(\mathbf{z})$ and whose covariance is a (fixed or learned) diagonal matrix:

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}\big(\mathbf{x};\, f_{\theta}(\mathbf{z}),\, \sigma_x^2 \mathbf{I}\big),$$

(2.5)

where $f_{\theta}(\mathbf{z})$ is the decoder's output (a deterministic nonlinear function of $\mathbf{z}$) and $\sigma_x^2$ is a variance term (often fixed or learned as a single scalar). In other contexts (e.g., image generation) $p_{\theta}(\mathbf{x}|\mathbf{z})$ can be modeled with a Bernoulli or other appropriate distribution. For discrete data, such as binary images, a Bernoulli distribution may be used:

$$p(x|z; \theta) = \prod_{i=1}^{D} p_i(z)^{x_i} (1 - p_i(z))^{1-x_i} \tag{2.6}$$

where $p_i(z)$ represents the probability of the $i$-th dimension being 1. The decoder aims to maximize the likelihood of the original data given $z$, ensuring that reconstructed outputs closely resemble the input, which is critical for applications like data generation and denoising .

During training, one first samples $z \sim q_{\phi}(z \mid x)$ from the encoder, then the decoder produces a reconstructed output $\hat{x}$ by sampling from $p_{\boldsymbol{\theta}}(x \mid z)$. The decoder thereby learns to invert the encoder's mapping, transforming latent codes back into the data space.

After training, the decoder network can be used on its own to gener-

ate new data: draw $z \sim \mathcal{N}(0, I)$ and then sample $x \sim p_{\boldsymbol{\theta}}(x \mid z)$. This generative ability is a key advantage of VAEs, allowing us, for example, to simulate plausible future price trajectories in a stock-market scenario by drawing many different $z$ samples [30].

### 2.2.3 Training and Sampling

Training a VAE involves learning both $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$ such that the model approximates the true data distribution and the encoder provides a good posterior approximation. The marginal likelihood, defined by Kingma and Max [22]

$$p_{\boldsymbol{\theta}}(x) = \int p_{\boldsymbol{\theta}}(x \mid z)\, p(z)\, dz, \qquad (2.7)$$

which is intractable for high-dimensional $z$. Instead, Using the fact that:

$$\log p_{\boldsymbol{\theta}}(x) = \log \frac{p_{\boldsymbol{\theta}}(x, z)}{q_{\boldsymbol{\phi}}(z|x)} + \log \frac{q_{\boldsymbol{\phi}}(z|x)}{p_{\boldsymbol{\theta}}(z|x)} \qquad (2.8)$$

which holds for any $z$, they maximize a tractable lower bound, the *evidence lower bound* (ELBO):

$$\log p_{\boldsymbol{\theta}}(x) = \underbrace{\mathbb{E}_{q_{\boldsymbol{\phi}}(z|x)}\big[\log p_{\boldsymbol{\theta}}(x \mid z)\big] - D_{\mathrm{KL}}\big(q_{\boldsymbol{\phi}}(z \mid x) \,\|\, p(z)\big)}_{\mathcal{L}(x; \boldsymbol{\theta}, \boldsymbol{\phi})}$$
$$+ D_{\mathrm{KL}}\big(q_{\boldsymbol{\phi}}(z \mid x) \,\|\, p_{\boldsymbol{\theta}}(z \mid x)\big). \qquad (2.9)$$

and since the KL term on the right is non-negative, $\mathcal{L}$ indeed lower-bounds $\log p_{\boldsymbol{\theta}}(x)$:

$$\mathcal{L}(x; \boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(z|x)}\big[\log p_{\boldsymbol{\theta}}(x \mid z)\big] - D_{\mathrm{KL}}\big(q_{\boldsymbol{\phi}}(z \mid x) \,\|\, p(z)\big). \qquad (2.10)$$

The ELBO comprises two competing objectives: (i) a *reconstruction term* that encourages the decoder to model the data accurately and (ii) a *regularization term* that pushes $q_\phi(z \mid x)$ toward the simple prior $p(z) = \mathcal{N}(0, I)$ [26]. Maximizing $\mathcal{L}$ thus teaches the decoder to generate realistic data while keeping the latent space well behaved.

In practice, the expected ELBO is maximized over the training set using stochastic gradient-based methods. Due to the additive nature of the ELBO across data points, stochastic mini-batch optimization is employed. For a given mini-batch, gradients of the ELBO are approximated and used to update the model parameters. The ELBO can be jointly optimized with respect to both generative parameters $\boldsymbol{\theta}$ and inference parameters $\boldsymbol{\phi}$ using stochastic gradient descent (SGD). For i.i.d. data, the ELBO over the dataset $\mathcal{D}$ decomposes as:

$$\mathcal{L}\boldsymbol{\theta}, \boldsymbol{\phi}(\mathcal{D}) = \sum x \in \mathcal{D}\mathcal{L}_{\boldsymbol{\theta},\boldsymbol{\phi}}(x). \tag{2.11}$$

Gradients with respect to $\boldsymbol{\theta}$ are straightforward to estimate using Monte Carlo sampling:

$$\nabla_{\boldsymbol{\theta}}\mathcal{L}\boldsymbol{\theta}, \boldsymbol{\phi}(x) \approx \nabla\boldsymbol{\theta} \log p_{\boldsymbol{\theta}}(x, z), \quad z \sim q_{\boldsymbol{\phi}}(z|x). \tag{2.12}$$

To compute gradients with respect to $\boldsymbol{\phi}$, the reparameterization trick [22] is used for continuous latent variables, allowing gradient flow through sampled latent vectors and obtaining unbiased, low-variance gradients by rewriting $z$ as:

$$z = \boldsymbol{\mu}\boldsymbol{\phi}(x) + \boldsymbol{\sigma}\boldsymbol{\phi}(x) \odot \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, I). \tag{2.13}$$

This results in a doubly stochastic optimization procedure—stochasticity arises both from the mini-batch sampling and the latent variable sampling, allows efficient end-to-end training of VAEs using standard back-propagation. This approach is commonly referred to as the Auto-Encoding Variational Bayes (AEVB) algorithm [22], [31]. Below is the pseudocode

for Auto-Encoding Variational Bayes (AEVB) algorithm introduced by Kingma et al [22]:

---

**Algorithm 1** Stochastic Optimization of the ELBO (Auto-Encoding Variational Bayes)

---

**Data:** Dataset $\mathcal{D}$, inference model $q_\phi(z|x)$, generative model $p_\theta(x, z)$
**Result:** Learned parameters $\theta$, $\phi$

  1: $(\theta, \phi) \leftarrow$ Initialize parameters
  2: **while** SGD not converged **do**
  3:       Sample minibatch $\mathcal{M} \subset \mathcal{D}$
  4:       Sample noise $\varepsilon \sim p(\varepsilon)$ for each $x \in \mathcal{M}$
  5:       Compute stochastic ELBO estimate $\hat{\mathcal{L}}_{\theta,\phi}(\mathcal{M}, \varepsilon)$
  6:       Compute gradients $\nabla_{\theta,\phi}\hat{\mathcal{L}}_{\theta,\phi}(\mathcal{M}, \varepsilon)$
  7:       Update $\theta$, $\phi$ using stochastic gradient descent

---

Gradients with respect to $\phi$ remain unbiased and low-variance thanks to the reparameterization trick introduced by Kingma and Max in 2013 [22], making this the standard *Auto-Encoding Variational Bayes* (AEVB) algorithm. The procedure converges when the model has learned both an effective generative distribution and a good inference mechanism.



Figure 2.5: Overview of a Variational Autoencoder (VAE). Image source: Understanding Deep Learning book

Once trained, the VAE can be used for *sampling* new data. Only the

decoder is required: sample $z \sim \mathcal{N}(0, I)$, then sample $x \sim p_{\boldsymbol{\theta}}(x \mid z)$. Repeating this yields plausible price paths, enabling scenario generation for risk assessment, trading-strategy back-testing, or dataset augmentation [24]. The combination of a probabilistic encoder, a generative decoder, and principled optimization via the ELBO (illustrated in Figure 2.5) is what enables VAEs to learn complex time-series distributions while keeping the latent space structured for generation [25], [26].

For models with discrete latent variables, the reparameterization trick is not directly applicable due to the non-differentiability of discrete sampling. Instead, an alternative unbiased gradient estimator known as the score function estimator [32] can be used to compute gradients of the ELBO. Given a function $f(z)$ and latent variable $z \sim q_\phi(z|x)$, the gradient with respect to $\boldsymbol{\phi}$ is given by:

$$\nabla_{\phi} \mathbb{E}{q\boldsymbol{\phi}(z|x)}[f(z)] = \mathbb{E}{q\boldsymbol{\phi}(z|x)}[f(z)\nabla_{\phi} \log q_\phi(z|x)], \qquad (2.14)$$

$$\approx f(z)\nabla_{\phi} \log q_\phi(z|x). \qquad (2.15)$$

This method—also known as the likelihood ratio estimator [33], [34] or REINFORCE [35]—has been widely adopted in settings involving discrete latents, including black-box variational inference [36], neural variational inference [37], and stochastic search variational Bayes [38]. Though often high in variance, this estimator remains broadly applicable and can be combined with control variates [39] for variance reduction.

# Chapter 3

# Literature Review

*This chapter reviews the current state-of-the-art in stock price prediction, with a focus on explainable AI (xAI). Specifically, I examine the latest advancements in VAEs models, considering their applications and contributions to the problem I am addressing.*

## 3.1   Stock price prediction

Stock price prediction represents a central problem in financial research, largely motivated by the potential for substantial economic gains. However, the inherent volatility, stochasticity, and non-linear nature of financial markets present significant modeling challenges. These complexities are further compounded by external factors, including macroeconomic events and public sentiment expressed through news and social media platforms. Over the past decades, a broad spectrum of models has been proposed to address these challenges, ranging from classical statistical frameworks to advanced deep learning architectures. Each modeling paradigm offers unique strengths, yet also faces distinct limitations with respect to capturing the intricate dynamics of financial time series.

Early approaches in this domain primarily relied on traditional time-series models, such as the Autoregressive Integrated Moving Average (ARIMA)

model [19], with the meaning explained as shown in Figure 3.1. ARIMA forecasts future stock values by modeling linear dependencies in historical price data under the assumption of stationarity. While ARIMA demonstrates satisfactory performance under stable market conditions, its linear structure limits its effectiveness in capturing abrupt market shifts and non-linear interactions, especially during periods of heightened volatility [1].



Figure 3.1: What is ARIMA model?. Image source: [40]

The advent of deep learning has facilitated the development of more expressive models capable of capturing complex temporal dependencies. Long Short-Term Memory (LSTM) networks, a variant of recurrent neural networks (RNNs), are particularly well-suited to sequential prediction tasks due to their ability to retain long-term contextual information. Fischer and Krauss [8] demonstrated that LSTM-based models significantly outperform traditional approaches in forecasting stock returns, highlighting their ability to extract meaningful temporal features from historical data. Despite these advances, LSTMs exhibit limitations in interpretability and

typically require engineered mechanisms to incorporate unstructured data such as textual sentiment.



(a) CSI 300

(b) S&P 500

(c) N225

(d) HSI

Figure 3.2: The result of experiments by Li et al. [41] demonstrate that Transformer outperforms other classic methods significantly.

To address the shortcomings of RNN-based methods, transformer architectures have been adapted for time-series forecasting, including stock price prediction tasks. Originally introduced in the context of natural language processing, transformers utilize self-attention mechanisms to capture both local and global dependencies in sequential data. Li et al. [41] proposed a transformer-based framework for predicting stock indices, leveraging multi-head attention and an encoder-decoder structure to model intricate market interactions, with results shown in Figure 3.2. These models often surpass LSTM baselines in terms of predictive accuracy and training efficiency

[42]. Nevertheless, standard transformer architectures are predominantly designed to process structured numerical data, and thus struggle to natively integrate unstructured textual inputs.

To capture the influence of external information sources, recent research has focused on models that jointly utilize historical price data and textual sentiment. StockNet, proposed by Xu and Cohen [16], employs a deep generative modeling framework that integrates Twitter data with stock price sequences through neural variational inference. The model introduces recurrent latent variables to account for uncertainty and temporal structure in financial data, achieving state-of-the-art results on movement prediction tasks by effectively modeling the interplay between textual sentiment and price dynamics.

In parallel, Hierarchical Attention Networks (HANs) have been introduced to extract salient information from textual data for stock prediction. Yang et al. [43] developed a HAN-based model that processes news articles at both the word and sentence level, using attention mechanisms to weigh relevant content. By fusing textual representations with technical indicators, the model enhances predictive accuracy and addresses challenges related to noisy or redundant text inputs.

Other hybrid methodologies incorporate textual features through classical or deep learning models. Schumaker and Chen [44] utilized Support Vector Machines (SVMs) in conjunction with features extracted from financial news, demonstrating that sentiment and event-driven information can materially influence price forecasting. Figure 3.3 showcases a simple results obtained through the application of their proposed method. Similarly, Ding et al. [45] proposed a neural tensor network that integrates structured event representations from news articles with price data to model event-driven market responses.

Despite these developments, several challenges persist in the integration of textual data into predictive frameworks. Textual information varies widely in quality and relevance, and extracting meaningful signals from

Figure 3.3: Example AZFinText representation by Schumaker and Chen et al. [44] .

unstructured sources remains non-trivial. Moreover, financial markets are highly dynamic; thus, models must be capable of rapid adaptation to newly emerging information—an aspect that remains an open problem even in state-of-the-art architectures

## 3.2 Explainable stock prediction

Explainable stock prediction has become a critical focus in financial forecasting, driven by the need for transparency in complex machine learning models to foster trust among investors, regulators, and analysts. The inherent opacity of deep learning models, such as Long Short-Term Memory (LSTM) networks and transformers, poses challenges in understanding the rationale behind predictions, which is essential for high-stakes financial applications. Explainable AI (XAI) techniques address this by pro-

viding insights into model decision-making processes, balancing predictive accuracy with interpretability. This subsection reviews key approaches to achieving explainability in stock prediction, focusing on attention-based methods, Local Interpretable Model-agnostic Explanations (LIME), SHapley Additive exPlanations (SHAP) and Large Language Model Application, each leveraging distinct mechanisms to enhance transparency while facing unique challenges.

Attention-based approaches [46]–[49], rooted in natural language processing, have been adapted to improve the interpretability of stock prediction models by highlighting the most influential input features or time steps. These mechanisms assign weights to input elements, such as historical prices or textual data from news and social media, indicating their relative importance in the prediction process. For instance, [46] introduced the Transformer Encoder-based Attention Network (TEANet), which integrates transformer models with attention mechanisms to extract deep features from small samples of financial data, including tweets and stock prices. By visualizing attention weights, TEANet reveals which textual or temporal elements drive predictions, achieving superior performance in stock movement forecasting compared to traditional methods. Similarly, [47] proposed the Attention Mechanism Variant LSTM (AMV-LSTM), which couples attention with LSTM to weigh the relevance of time steps, enhancing both accuracy and interpretability. However, attention-based models can be computationally intensive and may struggle with noisy textual inputs, necessitating robust preprocessing techniques [46].

Local Interpretable Model-agnostic Explanations (LIME) [51] provide a model-agnostic method for explaining individual predictions by approximating complex models with simpler, interpretable surrogates in the local vicinity of the prediction. Specifically, LIME generates perturbed samples around the instance of interest, queries the black-box model for predictions, and fits a locally weighted linear model to approximate the behavior of the original model near that instance, as shown in Figure 3.4. This

27

Figure 3.4: Steps of LIME algorithm. Image Source: [50]

process highlights the contribution of individual features to the prediction. In the context of stock prediction, LIME has been applied to models such as Random Forests to interpret the influence of features like closing prices, trading volumes, and technical indicators. For instance, [52] demonstrated that LIME effectively explains Random Forest predictions in stock price forecasting, identifying high trading volumes as significant predictors. However, while LIME offers valuable insights, its explanations can be sensitive to both the choice of the surrogate model and the definition of the local neighborhood, which may affect the consistency and reliability of the interpretations [51].

SHapley Additive exPlanations (SHAP) [53], rooted in cooperative game theory, offer a theoretically sound framework for interpreting model predictions. SHAP assigns each feature an importance score based on its contribution to the model's output relative to a baseline (typically the average prediction). The detailed calculation of this score (called Shapley value) is shown in Figure 3.5. This approach provides both local and

global interpretability, which is particularly beneficial in stock prediction tasks where understanding the broader relevance of features is crucial. For example, [54] utilized SHAP in conjunction with a deep learning model to forecast stock trends, identifying key indicators such as moving averages and volatility metrics as dominant contributors. Nevertheless, SHAP's computational complexity remains a limitation, especially when applied to large datasets or intricate models, thereby constraining its scalability for real-time financial applications [53].



Figure 3.5: Equation of Shapley value in an example. Image Source: [55]

More recently, Large Language Models (LLMs) have introduced a novel dimension to explainable stock prediction by generating natural language justifications for their outputs. Unlike traditional techniques that rely on numerical attribution scores or attention visualizations, LLMs are capable of producing coherent, human-readable explanations that articulate the reasoning behind predictions. For example, Koa et al [18] proposed the Summarize-Explain-Predict (SEP) framework, which leverages self-reflective agents and Proximal Policy Optimization to enable LLMs to

autonomously generate interpretable stock forecasts, as showcased in Figure 3.6. This method removes the need for manually curated explanation datasets while preserving high predictive accuracy.

**1. Task Prompt:**

```
Given a list of facts, estimate their overall impact on the price
movement of AAPL stock. Give your response in this format:
(1) Price Movement, which should be either Positive or Negative.
(2) Explanation, which should be in a single, short paragraph.

Facts: […]
```

**2. LLM Response:**

```
Price Movement: Mixed
Explanation: The given facts provide a mixed sentiment regarding
AAPL stock. On one hand, there are positive indicators such as the
stock being at an all-time high, Apple Watch outselling the Swiss
watch industry, and Apple's acquisition of NeXT. On the other
hand, there are negative aspects including the potential impact of
the coronavirus on AirPods orders, France's fine on Apple for
slowed iPhones, and the stock price falling after a downgrade.
Other facts provide neutral information about stock repurchases,
trademark appeals, and the growth of Google's Canadian workforce.
Overall, the sentiment is mixed, with both positive and negative
factors influencing the perception of AAPL stock.
```

Figure 3.6: Prediction and explaination generated by Summarize-Explain-Predict framework [18] from text prompts. Image source: [18]

Complementing this, Wang et al [56] introduced the LLMFactor framework, which employs Sequential Knowledge-Guided Prompting to extract economically meaningful factors from financial news and market data. Rather than relying solely on sentiment analysis, LLMFactor focuses on identifying latent drivers that directly influence stock movements. The framework has shown superior performance across several benchmark datasets from both U.S. and Chinese markets, underscoring the potential of LLMs to deliver both accurate and interpretable forecasts in complex financial

environments.

To further enhance interpretability, modern stock prediction systems increasingly incorporate hybrid architectures that combine multiple explainability techniques. These systems often integrate LSTM networks with attention mechanisms or transformer-based encoders to offer temporal and feature-wise attributions. The CLEAR-Trade framework [57] exemplifies this approach, providing interactive visualization and explanation tools while maintaining a prediction accuracy of 61.22

Moreover, many of these systems employ multi-modal data inputs, merging historical price series, technical indicators, sentiment scores, and macroeconomic variables to construct richer, more explainable models. The inclusion of explainability constraints during training ensures that interpretation mechanisms are not post-hoc add-ons but are instead embedded into the model's optimization process, thereby producing more faithful and informative explanations.

In summary, interpretability methods such as attention mechanisms, LIME, SHAP, LLMs, and hybrid models have significantly advanced the transparency of stock prediction systems. Attention-based techniques help uncover temporal and contextual relevance, while LIME and SHAP provide flexible, model-agnostic insights into local and global feature importance. LLMs add a powerful narrative layer to explanations, and hybrid frameworks unify these approaches for comprehensive interpretability. Despite these advancements, challenges such as computational demands, sensitivity to input noise, and stability of explanations persist, highlighting the ongoing need for robust, efficient, and scalable explainable AI (XAI) solutions in the financial domain.

While existing approaches to explainable stock prediction—such as LIME, SHAP, attention-based models, large language models, and hybrid architectures—have made significant strides in enhancing model transparency, they largely rely on correlative reasoning and fail to explicitly distinguish between causally predictive signals and reactive noise. Notably,

current literature lacks methods that integrate causal attention mechanisms with volatility-aware representations, particularly within the context of explainable financial forecasting. This gap limits the interpretability and reliability of predictions in highly dynamic and uncertain market conditions.

To address this limitation, in this thesis, I propose the Causal and Dual-Pathway Prediction-Explanation Network (CDP-PEN), a novel architecture designed to simultaneously improve predictive accuracy and interpretability in text-based stock forecasting based on the PEN architecture, with the overall architecture mentioned in section 3.3 proposed by Li et al [17]. The model introduces two core components: (1) the Causal Text Selection Unit (Ca-TSU), which filters textual embeddings to isolate those with causal influence, producing a Vector of Causality (VoC) that quantifies the contribution of each input feature; and (2) the Dual-Pathway Shared Representation Learning (DP-SRL) module, which explicitly models the bidirectional interaction between text and market behavior. DP-SRL incorporates pathway-specific attention mechanisms to disentangle predictive signals from reactive patterns, gated inter-pathway communication to regulate information flow, and volatility-aware fusion gates to dynamically adapt feature weighting based on market uncertainty. By integrating causality and volatility sensitivity into a unified explainable framework, CDP-PEN aims to advance the frontier of interpretable and robust financial time-series modeling.

## 3.3   Overview of PEN

### 3.3.1   Overall architecture

The Prediction-Explanation Network (PEN), introduced by Li et al. [17], is a novel framework designed to predict stock price movements while enhancing explainability by aligning textual data (e.g., news articles, so-

Figure 3.7: The overall framework of PEN model. Image source: [17]

cial media posts) with historical price sequences. Figrure 3.7 shows the overall architecture of PEN. Addressing the challenge of modeling complex market dynamics driven by multimodal inputs, PEN integrates text and price data through a shared representation learning approach, achieving state-of-the-art accuracy and interpretability. Its primary goal is to deliver accurate predictions of stock price movements (binary classification: upward or downward) while providing transparent explanations via a Vector of Salience (VoS), which quantifies the importance of individual text documents. This section re-explains PEN's main components—Text Embedding Layer (TEL), Shared Representation Learning (SRL) with its subcomponents, Deep Recurrent Generation (DRG), and Temporal Attention Prediction (TAP)—and highlights its correlation-based limitation, set-

33

ting the stage for advancements proposed in the Causal and Dual-Pathway Prediction-Explanation Network (CDP-PEN).

**Text Embedding Layer (TEL)**

The Text Embedding Layer (TEL) processes textual inputs to generate lower-dimensional representations that capture contextual information. For a given stock on trading day $t$, the text corpus $\mathcal{C}_t = \{C_{t1}, C_{t2}, \ldots, C_{tM}\} \in \mathbb{R}^{M \times l \times h_w}$ consists of $M$ documents, where each document $C_{tM}$ is a sequence of $l$ words with word embedding size $h_w$. PEN employs a bidirectional Gated Recurrent Unit (Bi-GRU) to encode each document, capturing both past and future word contexts. The output is a matrix of text embeddings $e_t \in \mathbb{R}^{h \times M}$. The Bi-GRU ensures that semantic and contextual nuances in texts, such as sentiment or event-driven information, are effectively represented for downstream processing.

**Shared Representation Learning (SRL)**



Figure 3.8: The detailed structure of SRL module. Authors use information, including texts and prices, of 02/07- 04/07 to predict stock movement of 05/07 for illustration. Image source: [17]

The Shared Representation Learning (SRL) module, as shown in Fig-

ure 3.8, is the cornerstone of PEN, modeling the interaction between text embeddings $e_t$ and price data $p_t \in \mathbb{R}^{L \times 3}$, where $p_t$ includes normalized adjusted close, high, and low prices over a lag window $L$. SRL generates a shared representation $h_t \in \mathbb{R}^{h \times 1}$ and a Vector of Salience (VoS) $\omega_t \in \mathbb{R}^{M \times 1}$, which quantifies the importance of each text document for prediction. The original PEN architecture incorporates a Shared Representation Learning (SRL) module that aligns text streams and price streams via a salience-based attention mechanism. While effective in selecting relevant messages, it is correlation-based and lacks causality awareness, limiting its ability to distinguish predictive from reactive texts. SRL comprises three subcomponents:

- **Text Selection Unit (TSU)**: The TSU filters text embeddings to identify those most relevant to stock price movements, mitigating noise from irrelevant or low-quality texts. It computes the VoS using a softmax over a weighted combination of the previous hidden state $h_{t-1}$ and text embeddings $e_t$:

$$\omega_t = \mathrm{softmax}\left(k_t^\top \tanh\left(W_1 h_{t-1} + W_2 e_t + b_1\right)\right) \qquad (3.1)$$

  The weighted text embedding $i_t = e_t \omega_t$ emphasizes salient texts, enhancing prediction accuracy and explainability.

- **Text Memory Unit (TMU)**: Inspired by LSTM cells, the TMU preserves temporal information from text embeddings, recognizing that market responses to news may be delayed. It uses forget and output gates to regulate information flow:

$$
\begin{aligned}
f_t &= \sigma\left(W_3[i_t, h_{t-1}] + b_3\right) & o_t &= \sigma\left(W_4[i_t, h_{t-1}] + b_4\right) \\
l_t &= \tanh\left(W_5[i_t, h_{t-1}] + b_5\right) & l_t &= f_t \odot l_{t-1} + o_t \odot l_t & (3.2)
\end{aligned}
$$

  Here, $l_t \in \mathbb{R}^{h \times 1}$ is the text memory state, initialized via the Xavier

algorithm, ensuring retention of relevant historical text information.

- **Information Fusion Unit (IFU)**: The IFU integrates text and price data into a shared representation $h_t$. It employs a gating mechanism to balance contributions:

$$
\begin{aligned}
d_t &= \sigma \left( W_6[p_t, l_t, h_{t-1}] \right) & h_p &= \tanh \left( W_8[p_t, h_{t-1}] \right) \\
h_l &= \tanh \left( W_7[l_t, h_{t-1}] \right) & h_t &= d_t \odot h_p + (1 - d_t) \odot h_l
\end{aligned} \tag{3.3}
$$

The resulting $h_t$ is passed to the next SRL module, forming a recurrent structure that captures temporal dependencies across days. Authors regard the last hidden state $h_t$, i.e. the hidden state of the last day in time lag, as input of deep recurrent generation module.

## Deep Recurrent Generation (DRG)

The Deep Recurrent Generation (DRG) module employs a variational autoencoder (VAE) to model stock price movements probabilistically, capturing stochastic market factors. It takes the shared representation $h_t$ from SRL as input $x_t$ and infers a latent variable $z_t$. The VAE approximates the conditional probability $p_\theta(y \mid X)$ using a recurrent structure with Gaussian-distributed prior and posterior:

$$
p_\theta(z_t \mid z_{<t}, x_{\leq t}) \sim \mathcal{N}(z_t; \mu_\theta, \sigma_\theta^2 I), \quad q_\phi(z_t \mid z_{<t}, x_{\leq t}, y_t) \sim \mathcal{N}(z_t; \mu_\phi, \sigma_\phi^2 I) \tag{3.4}
$$

A GRU-based encoder processes $x_t$ to produce hidden states $h_t^{\text{enc}}$, which are used to compute the mean and variance of the posterior and prior distributions via reparameterization. The decoder generates predictions $\hat{y}_t$ (for $t < T$) by combining the latent variables and shared representations, enabling robust modeling of market uncertainty.

**Temporal Attention Prediction (TAP)**

The Temporal Attention Prediction (TAP) module refines predictions by leveraging temporal dependencies across previous predictions. It uses the decoder hidden states $H^{\text{dec}} = [h_1^{\text{dec}}, \ldots, h_{T-1}^{\text{dec}}]$ to compute attention weights:

$$q^{\text{dec}} = (h_{T-1}^{\text{dec}})^\top \tanh(W_q H^{\text{dec}}) \qquad k^{\text{dec}} = w_k^\top \tanh(W_k H^{\text{dec}})$$

$$w^{\text{dec}} = \text{softmax}\left(q^{\text{dec}} \odot (k^{\text{dec}})^\top\right) \qquad v^{\text{dec}} = [\hat{y}_1, \hat{y}_2, \ldots, \hat{y}_{T-1}] \qquad (3.5)$$

The final prediction $\hat{y}_T$ is generated by applying a softmax over a weighted combination of past predictions and the current hidden state:

$$\hat{y}_T = \text{softmax}\left(W_{\text{dec}}^\top [v^{\text{dec}}(w^{\text{dec}})^\top, h_T^{\text{dec}}] + b_{\text{dec}}\right) \qquad (3.6)$$

This mechanism ensures that the model captures temporal correlations, enhancing prediction accuracy.

## 3.3.2 Goal of PEN

PEN's dual objective is to achieve high predictive accuracy and provide interpretable explanations for stock price movements. By generating the VoS through the SRL module, PEN identifies the most relevant texts influencing predictions, aligning with investor decision-making processes that focus on key information. Experimental results on the ACL18 and DJIA datasets demonstrate that PEN outperforms baselines like Stock-Net, achieving accuracies of 59.90% and 60.51%, respectively, while the VoS achieves a 99.5% Ratio of Top Two (RTT) and 89.3% Ratio of Relevance (RoR), far surpassing attention-based methods [17]. However, the correlation-based nature of the SRL module limits its ability to distinguish causal from reactive texts, a gap addressed by the proposed CDP-PEN, which introduces causal reasoning and volatility-aware mechanisms to en-

hance both accuracy and interpretability.

### 3.3.3 Limitations of PEN and Motivation for Improvement

The Prediction-Explanation Network (PEN) represents a significant advancement in stock price movement prediction by integrating text and price data through its Shared Representation Learning (SRL) module, achieving superior accuracy and explainability via the Vector of Salience (VoS) [17]. However, several limitations in its design restrict its robustness and interpretability, motivating the development of the Causal and Dual-Pathway Prediction-Explanation Network (CDP-PEN).

A primary limitation of PEN is that the VoS, generated by the SRL's Text Selection Unit (TSU), relies on correlation-based attention mechanisms to identify salient texts. This approach is susceptible to noise in textual data, such as irrelevant or low-quality news and social media posts, which can skew predictions and explanations. Although PEN effectively highlights salient messages, it does not differentiate between causally predictive texts and merely reflective ones, leading to potential misattribution of influence. For instance, reactive texts describing past price movements may be weighted similarly to predictive texts signaling future trends, reducing the model's explanatory clarity.

Furthermore, the SRL module does not explicitly distinguish between predictive and reflective messages, resulting in unclear explanations about which texts drive price movements and which merely describe them. This lack of causal reasoning limits PEN's ability to provide transparent, actionable insights for investors and regulators, who require precise attribution to meet compliance standards, such as the EU's MiFID II.

Additionally, PEN lacks mechanisms to adapt to volatile market conditions. The SRL and Temporal Attention Prediction (TAP) modules do not account for varying market regimes (e.g., bull, bear, or high-volatility

38

periods), limiting the robustness of its explanatory power and predictive accuracy in dynamic environments. This is particularly critical in financial markets, where sudden shifts triggered by external events demand adaptive modeling.

These limitations motivate the development of CDP-PEN, which introduces a Causal Text Selection Unit (Ca-TSU), with details are covered in Chapter 4 to filter texts based on causal influence, producing a Vector of Causality (VoC) to enhance interpretability. Additionally, CDP-PEN's Dual-Pathway Shared Representation Learning (DP-SRL), mentioned in Chapter 5 explicitly separates predictive and reflective texts, while a Volatility-Aware Fusion Mechanism, introduced in Chapter 6 dynamically adjusts input weighting based on market conditions. These improvements aim to address PEN's shortcomings, offering a more robust, interpretable, and adaptable framework for stock price prediction.

# Chapter 4

# Causal Text Selection Unit (Ca-TSU)

*This chapter introduces the Causal Text Selection Unit (Ca-TSU), the first major contribution of this thesis. This module filters textual inputs, such as news and social media, to isolate messages with causal influence on stock prices. By employing a specialized attention mechanism, the Ca-TSU prioritizes predictive signals, reducing the impact of irrelevant or noisy text. The chapter details the implementation within the proposal framework, including the configuration of attention weights and regularization strategies. Experimental results demonstrate the effectiveness of the Ca- TSU in enhancing prediction accuracy by focusing on high-impact textual data, thereby addressing the challenge of noise in unstructured inputs.*

## 4.1 Overview

In this chapter, we introduce the **Causal Text Selection Unit (Ca-TSU)**, a novel component designed to enhance the interpretability and effectiveness of text data usage in stock price prediction. The Ca-TSU extends the original Text Selection Unit (TSU) of the Prediction-Explanation

Network (PEN) [17] by incorporating principles of causal inference into the text selection process. Traditional attention-based mechanisms in stock prediction models often identify texts that correlate with future price movements [58], [59], but they may not distinguish whether those texts have a genuine *causal* impact on the price or are merely coincidental. Our proposed Ca-TSU addresses this limitation by explicitly aiming to select news or social media texts that are not only relevant but also *causally* influential for price movements. By doing so, Ca-TSU filters out noisy or irrelevant texts and focuses on a subset of information that can truly explain and predict stock movements, aligning with the intuition that investors focus on key causal factors rather than all available information.

From an academic standpoint, Ca-TSU is motivated by the need for more interpretable AI in finance and the concept that explainability can be improved by aligning model attention with causal effects. Prior studies have highlighted that not all texts contribute equally to market movements and that selecting a small subset of crucial texts can improve both predictive performance and human interpretability [5], [60]. Ca-TSU builds on this insight by integrating a causal selection strategy into the model's architecture. In summary, the Ca-TSU is an essential innovation for making the stock prediction model more robust to textual noise and more aligned with cause-effect relationships present in financial news, which we detail in the following sections.

## 4.2 Architecture

The architecture of the Ca-TSU module is illustrated in Figure 4.1. At a high level, Ca-TSU operates within the shared representation learning framework of my stock prediction model, replacing the standard TSU [17] with an enhanced unit that emphasizes causal relevancy. The Ca-TSU takes as input a set of text embeddings $\{\mathbf{e}_{t,1}, \mathbf{e}_{t,2}, \ldots, \mathbf{e}_{t,M}\}$ for all $M$ news articles or tweets on day $t$, as well as the previous day's hidden state $h_{t-1}$

from the model's recurrent pipeline (which encodes prior information from both price and text modalities). Internally, Ca-TSU uses an attention mechanism similar to the original TSU to compute a *Vector of Salience* (VoS) $\omega_t \in \mathbb{R}^M$ that indicates the importance of each candidate text on day $t$. For each time step $t$ and message $i$, the model computes a content-based attention score as follows:

$$a_{t,i} = \tanh\Big(W_1\, h_{t-1} \;+\; W_2\, \mathbf{e}_{t,i} \;+\; b_1\Big) \tag{4.1}$$

The attention weights over messages are then computed as:

$$\omega_{t,i} = \mathrm{softmax}\,(k_t a_{t,i}) \tag{4.2}$$

where $\mathbf{e}_{t,i} \in \mathbb{R}^d$ is the embedding of the $i$-th text (of dimension $d$), and $W_1, W_2$ are weight matrices that project the previous hidden state $h_{t-1}$ and the text embedding $\mathbf{e}_{t,i}$ into a common space. The vector $k_t$ and bias $b_1$ are learnable parameters that together determine the relevance of each text given the context. However, unlike a standard attention unit, Ca-TSU's attention computation is augmented with a causal perspective. This is achieved by incorporating features that reflect the potential impact of each text on future price changes. For example, Ca-TSU can include an additional input signal reflecting recent price volatility or trend direction, ensuring that texts are scored highly not just because they are contextually relevant to $h_{t-1}$, but because they are likely to *drive* upcoming price movement.

Concretely, Ca-TSU computes an importance score $\alpha_{t,i}$ for each text $i$ on day $t$ using a neural attention function that considers both the textual content and the prior state. To capture causal dependencies among messages, Ca-TSU introduces a causal attention mechanism:

$$q_{t,i} = e_{t,i} W_Q \tag{4.3}$$

$$k_{t,i} = e_{t,i} W_K \tag{4.4}$$

where $W_Q$ and $W_K$ are learnable matrices that project the message embeddings into query and key spaces, respectively.

The causal attention score is then computed as:

$$\alpha_{t,i} = \text{softmax}\left(q_{t,i}^{\top} \tanh(k_{t,i} + b_{causal})\right) \tag{4.5}$$

where $b_{causal}$ is a bias term. The inner product $q_{t,i}^{\top} \tanh(k_{i,t} + b_{causal})$ measures the compatibility between the query and the (biased) key, and the softmax normalizes these scores across messages. By using the same message embedding for both the query and key, but allowing for different learned projections and a bias, the model can capture not only the content relevance but also potential causal relationships among messages. For example, certain messages may be more influential not just because of their content, but because of how they interact with the latent structure of the message set on that day.

After computing raw scores $\alpha_{t,i}$ for all texts, Ca-TSU combines the content-based and causal attention mechanisms, the final attention weights are computed and normalized as:

$$\omega_{t,i}^{\prime \text{adj}} = \frac{\omega_{t,i} \cdot \alpha_{t,i}}{\sum_{i'} P_{t,i'} \cdot \alpha_{t,i'} + \epsilon} \tag{4.6}$$

where $\epsilon$ is a small constant for numerical stability. This fusion allows the model to attend to messages that are both content-relevant and causally important. The resulting vector $P_t = [\omega_{t,1}', \omega_{t,2}', \ldots, \omega_{t,M}']^{\top}$, which can be called *Vector of Causality* (VoC) reflects a probability distribution over the texts, effectively selecting the most important pieces of information on day $t$. A key difference in Ca-TSU is that this distribution is expected to be *sharply peaked* on truly influential news. In other words, Ca-TSU is architected to highlight one or a few documents that have the highest causal significance, rather than spreading attention widely across many

texts. Once the salience weights are obtained, we compute a **causal text representation** $i_t$ for the day by aggregating text embeddings:

$$i_t = \sum_{i=1}^{M} \omega_{t,i}'^{\text{adj}} \mathbf{e}_{t,i}. \tag{4.7}$$

This $i_t \in \mathbb{R}^d$ is a single vector summarizing the day's selected information. It emphasizes content from the most salient (and ideally causal) text messages, effectively filtering out the less relevant news. The vector $i_t$ then flows into the next components of the model (the Text Memory Unit and the fusion mechanism of the shared representation module) as the representative textual input for day $t$, and then this representative textual input is process through the variant of Variational Auto-Encoders (VAE) to generate stock movements from latent variables as Li at al. [17] proposed.

Architecturally, Ca-TSU is implemented as a dedicated module that can be plugged into our model's pipeline in place of a standard attention unit. The design ensures modularity: Ca-TSU can be toggled on or off, which is useful for ablation studies. In summary, the Ca-TSU architecture refines the model's focus on text data by combining attention with causal insight, yielding a more interpretable and noise-resistant input for prediction.

## 4.3 Learning Objective

The learning objective associated with the Ca-TSU module aligns with the overarching goals of the prediction-explanation framework: we seek to maximize prediction accuracy while simultaneously maximizing the explainability of the model's decisions. Ca-TSU contributes to both of these goals. In terms of prediction accuracy, the causal filtering of texts ensures that the model is fed with more informative signals, which should improve its ability to predict stock movements. In terms of explainability, Ca-TSU's output—the Vector of Causality $\omega_t'$ highlighting a few key

texts—serves as a direct explanation for the model's prediction, identifying which news items are considered causes for the predicted price movement.

Formally, let $\mathcal{L}_{\text{pred}}$ denote the primary prediction loss (the variational lower bound-ELBO in the case of a VAE-based predictor as described in Chapter 2 and [16], [17]):

$$\mathcal{L}_{\text{pred}} = \mathbb{E}_{q(z|x,y)}\left[\log p(y|g)\right] - \lambda_{\text{KL}} \cdot \text{KL}\left(q(z|x,y)\|p(z|x)\right) \tag{4.8}$$

where $q(z|x, y)$ is the approximate posterior over latent variables $z$ given inputs $x$ and targets $y$, $p(y|g)$ is the likelihood of the target given the latent representation $g$ and $\lambda_{\text{KL}}$ is a coefficient (possibly annealed during training) that controls the strength of the KL regularization. The first term encourages the model to explain the observed data, while the second term regularizes the latent space.

I then add an explainability-oriented regularization term $\mathcal{L}_{\text{consistency}}$ to the loss to encourage the desired behavior in Ca-TSU's attention output. Following the approach in the original PEN model [17], I define $\mathcal{L}_{\text{consistency}}$ to promote that the content-based and causal attention mechanisms produce *consistent distributions*, by using a Jensen-Shannon (JS) divergence:

$$M = 0.5(\omega + \alpha) \tag{4.9}$$

$$\mathcal{L}_{\text{consistency}} = \frac{1}{2}\left[\text{KL}(\omega\|M) + \text{KL}(\alpha\|M)\right] \tag{4.10}$$

where:

- $\omega$ and $\alpha$ are the (clipped) content and causal attention weights, respectively.

- $M$ is the average of the two distributions.

- The JS divergence penalizes discrepancies between the two attention mechanisms, encouraging them to agree on which messages are important.

To prevent the attention distributions from becoming overly sharp (i.e., focusing too much on a single message), entropy regularization is applied:

$$\mathcal{L}_{\text{entropy}} = \mathbb{H}[P] + \mathbb{H}[\alpha] \tag{4.11}$$

where $\mathbb{H}[P]$ and $\mathbb{H}[\alpha]$ denote the entropy of the content and causal attention distributions, respectively. This term encourages the model to maintain a certain level of uncertainty in its attention, which can improve generalization and interpretability.

The overall training objective for Ca-TSU is formulated as a joint loss that balances predictive accuracy with the quality and consistency of the model's attention over textual inputs. Specifically, the total loss is given by:

$$\mathcal{L} = \mathcal{L}_{\text{pred}} + \lambda_{\text{consistency}} \, \mathcal{L}_{\text{consistency}} + \beta \, \mathcal{L}_{\text{entropy}}, \tag{4.12}$$

where $\mathcal{L}_{\text{pred}}$ is the main prediction loss (e.g., negative log-likelihood or variational evidence lower bound), $\mathcal{L}_{\text{consistency}}$ is a regularization term that encourages agreement between the content-based and causal attention distributions, and $\mathcal{L}_{\text{entropy}}$ is an entropy regularizer that promotes sparsity and interpretability in the attention weights. The hyperparameters $\lambda_{\text{consistency}}$ and $\beta$ control the trade-off between prediction accuracy, attention consistency, and explanation sharpness.

A proper choice of these hyperparameters is crucial: if $\lambda_{\text{consistency}}$ or $\beta$ are set too high, the model may overemphasize attention regularization at the expense of predictive performance; if too low, the model may ignore the regularizers and behave like a standard attention model. In my experiments, I tune these coefficients to ensure that Ca-TSU produces interpretable, causally meaningful attention distributions without sacrificing forecasting accuracy.

Importantly, Ca-TSU does not require explicit supervision regarding which texts are truly causal. Instead, the model learns to assign higher attention weights to texts that consistently contribute to reducing the pre-

diction loss $\mathcal{L}_{\mathrm{pred}}$. The consistency regularizer $\mathcal{L}_{\mathrm{consistency}}$ further aligns the content-based and causal attention mechanisms, while the entropy regularizer $\mathcal{L}_{\mathrm{entropy}}$ encourages the model to focus its attention on the most informative messages and suppress irrelevant ones. As a result, Ca-TSU naturally identifies causally relevant texts by optimizing for both predictive accuracy and interpretable, concentrated attention.

This approach is related to the concept of "causal attention" in recent literature [61], where attention mechanisms are guided to reflect underlying causal relationships. In Ca-TSU, this principle is embedded directly into the objective function through the use of consistency and entropy-based regularization, providing a simple yet effective means of promoting causal interpretability in text-based time series forecasting.

In summary, the Ca-TSU learning objective jointly optimizes for accurate prediction and interpretable, causally meaningful text selection. By incorporating attention consistency and entropy regularization, Ca-TSU steers the model to leverage textual information in a way that is both effective for forecasting and transparent for explanation.

## 4.4 Experiments

### 4.4.1 Experimental Setup

**Datasets**

I evaluate the effectiveness of our Ca-TSU model on two publicly available datasets: the ACL18 dataset [16] and the Daily News for Stock Price Movement Prediction Dataset (DJIA)[1]. I select these datasets for the following reasons: (1) they span distinct time periods, providing temporal diversity in our evaluation; (2) ACL18 focuses on individual stocks while DJIA targets stock market indices, offering different levels of market granu-

---

[1]https://www.kaggle.com/aaron7sun/stocknews

larity; and (3) they incorporate two fundamentally different types of textual data: social media posts (tweets) and news articles, respectively.

The ACL18 dataset comprises text data and historical price information for 88 highly traded US stocks across 9 industries, covering the period from 2014-01-01 to 2016-01-01. The textual content consists of tweets retrieved from Twitter, while the historical price data are sourced from Yahoo Finance. We process the ACL18 dataset following the methodology proposed in (Xu and Cohen 2018). The DJIA dataset contains news articles and price data for the Dow Jones Industrial Average from 2008-06-08 to 2016-07-01, where the news data consists of the top 25 headlines from the Reddit WorldNews Channel for each trading day.

## Evaluation Metrics

Following the evaluation protocol established by Xu and Cohen [16], I assess model performance using two primary metrics: accuracy (ACC) and Matthews Correlation Coefficient (MCC). Accuracy measures the proportion of correct predictions, while MCC provides a balanced measure that accounts for class imbalance and is particularly suitable for binary classification tasks in financial prediction scenarios.

## Parameter Settings

I implement our Ca-TSU model using TensorFlow to construct the computational graph. All bias terms are initialized to zero, while weights are initialized using the Xavier algorithm (Glorot and Bengio 2010). The model is optimized using the Adam optimizer with a learning rate of 1e-3. We employ a batch size of 32 with shuffled samples and utilize a 5-day lag window to enable the model to learn from historical context. The hidden state sizes in the Ca-TSU module and word embedding layer are set to 100 and 50 dimensions, respectively. To regularize the latent variables and prevent overfitting, we apply an input dropout rate of 0.3. Given the relatively

small size of the DJIA dataset, we adopt a transfer learning approach by pre-training models on the ACL18 dataset and then fine-tuning them on the DJIA dataset for all baseline comparisons.

**Baselines**

I compare our Ca-TSU model against the following established baselines:

- **Random:** A baseline that generates random movement predictions.

- **Random Forest (RF) (Pagolu et al. [62]):** A Random Forest classifier that incorporates sentiment analysis with Word2Vec embeddings.

- **HAN (Hu et al. [58]):** A hybrid attention network that leverages related news articles for prediction.

- **Stocknet (Xu and Cohen [16]):** A deep generative model that jointly exploits textual and price signals.

- **CPC (Wang et al. [63]):** A copula-based contrastive predictive coding method that models relevant macroeconomic variables to enhance prediction accuracy.

- **PEN (Li et al. [17]):** Prediction-Explanation Network to Forecast Stock Price Movement with Better Explainability

These baselines represent a diverse range of approaches, from traditional machine learning methods to state-of-the-art deep learning architectures, providing a comprehensive evaluation framework for my proposed Ca-TSU model.

I evaluate the impact of the Causal Text Selection Unit through a series of experiments, including ablation studies and performance comparisons on real-world stock movement prediction datasets. Our experiments aim to

answer two key questions: (1) Does Ca-TSU improve prediction accuracy compared to the baseline model without Ca-TSU? (2) Does Ca-TSU improve the explainability of the model's predictions by selecting better (more causal) texts?

**Prediction Accuracy:** I first measure the model's stock movement prediction accuracy with and without the Ca-TSU module. In these experiments, all other components of the model are kept identical, isolating the effect of Ca-TSU. We observe that integrating Ca-TSU yields approximately a **2% improvement in accuracy** over the baseline PEN model. As shown in Table 4.1, Ca-TSU consistently outperforms all baselines, including the state-of-the-art PEN model, across all evaluation metrics.. This improvement of roughly 2 percentage points is consistent across multiple datasets and market scenarios. The boost in performance can be attributed to Ca-TSU's ability to discard distracting or irrelevant texts (which may have acted as noise for the baseline) and focus on truly informative content that drives stock movements.

Table 4.1: Performance Comparison of Models on ACL18 and DJIA Datasets

| Models | ACL18 | | DJIA | |
|--------|--------|--------|--------|--------|
| | ACC(%) | MCC | ACC(%) | MCC |
| Random | 48.7600 | -0.002142 | 49.2245 | 0.0003192 |
| RF | 50.0430 | 0.010023 | 51.3266 | 0.050091 |
| HAN | 54.3270 | 0.052312 | 51.3409 | 0.059321 |
| StockNet | 54.5490 | 0.079726 | 52.9091 | 0.129039 |
| CPC | 54.6120 | **0.178910** | - | - |
| PEN | 54.9876 | 0.153219 | 55.5327 | 0.220024 |
| **Ca-TSU (Mine)** | 57.4405 | 0.162312 | **58.1367** | **0.234681** |

**Ablation Study on Text Selection:** To further validate Ca-TSU's contribution, I conduct an ablation study where we compare the full model against a variant where Ca-TSU is replaced by a standard TSU (or a simple attention mechanism) lacking the causal enhancement. The results confirm that the full model (with Ca-TSU) outperforms the variant in both accuracy and stability of predictions. Moreover, we notice that the variant without Ca-TSU is more prone to overfitting to spurious textual patterns—e.g., with an example shown in Figure 4.2 it might latch onto frequently occurring words or news that correlate with the target in the training set but are not truly predictive in general. In contrast, Ca-TSU's focus on causal relevance acts as a form of regularization, guiding the model to rely on more robust text signals.

**Explainability and Selected Texts:** A crucial part of our experimental evaluation involves assessing the explainability of model predictions. I use metrics proposed in [17] to quantify how well the model's selected texts align with human intuition and known market drivers. These metrics include:

- **RTT (Relevant Text Proportion)**: the proportion of top-ranked texts that are deemed relevant by human evaluators or by retrospective analysis.

- **RoR (Rate of Return of explanation)**: an information-retrieval inspired metric that considers how quickly the model's ranked list of texts "finds" the truly relevant news among all available news.

- **Kappa Agreement ($\kappa$)**: the inter-rater agreement between the model's selected texts and human experts, measured by Cohen's or Fleiss' $\kappa$ [17] over many instances.

In my experiments, Ca-TSU significantly improves these explainability metrics compared to baseline attention. For instance, we find that in

many cases Ca-TSU's top-weighted text for a given stock movement coincides with the news article that financial analysts later identify as the likely cause of that movement (e.g., a surprise earnings announcement or a major geopolitical event). The $\kappa$ score between Ca-TSU's selections and three human annotators' picks increases by a notable margin, indicating better agreement on what information matters. Additionally, Ca-TSU often achieves a high RTT with just 1 or 2 texts, whereas the baseline might need 3-5 texts to cover the key information.

Quantitatively, on ACL18 dataset I evaluated, the model with Ca-TSU achieved a top-1 relevant text accuracy of 85%, meaning in 85% of test instances the single highest-weight news was deemed the correct causal news by experts, compared to 75-80% for the baseline model. Such results demonstrate Ca-TSU's effectiveness in not only predicting the correct movement but also in providing a concise explanation for that prediction.

**Robustness Tests:** I also test the robustness of Ca-TSU under various conditions. One concern might be whether Ca-TSU could miss relevant information if there are multiple independent causes for a price move (e.g., two different news events affecting the stock simultaneously). In these scenarios, I found that Ca-TSU can sometimes split attention between two texts if needed (the entropy regularizer does not force a single text if two texts are truly important). The model is still able to handle multi-cause situations, though it naturally leans towards highlighting the dominant cause. We also experimented with different levels of news noise (adding irrelevant articles) and found that Ca-TSU scales well: even as we increase the number of distracting news articles, Ca-TSU continues to pick out the correct ones with high precision, whereas a non-causal attention baseline starts to degrade in performance by spreading weights among many irrelevant inputs.

Overall, the experiments confirm that Ca-TSU delivers on its design goals. It provides about a 2% absolute accuracy improvement over the

baseline, and it substantially enhances the clarity and correctness of the model's text-based explanations. These benefits come without a significant increase in model complexity—Ca-TSU adds only a small number of parameters (for the additional weight matrices like $W_2$ and possibly causal feature weights) and negligible computational overhead compared to the overall model.

## 4.5 Applications

The Causal Text Selection Unit has several practical applications in the realm of financial prediction and beyond. By improving both the predictive power and transparency of the model, Ca-TSU enables a range of use cases:

- **Real-time News Filtering for Traders:** In live trading systems, there is an overwhelming flow of news. Ca-TSU can serve as an intelligent filter that flags the news most likely to impact stock prices. Traders and analysts can use the model's selected headlines as a quick situational awareness tool, focusing their attention on potentially market-moving information and ignoring background noise.

- **Model Explainability for Regulatory Compliance:** Financial institutions often face regulatory requirements to explain AI-driven decisions (for example, why a trading algorithm made a certain buy/sell decision). With Ca-TSU, whenever the model predicts a significant price movement or recommends an investment action, it also provides the top causal text(s) behind this decision. Such capability can be crucial for compliance and auditing, as it offers a traceable justification grounded in news events or social media discussions, aligning model decisions with real-world information.

- **Decision Support for Investors:** Beyond automated trading, Ca-TSU can be used in decision support tools for human portfolio man-

agers. It can highlight key pieces of news driving a stock's momentum. For instance, an investor using a portfolio analysis platform might get explanations like: "Stock X is predicted to rise because of news Y (selected by Ca-TSU)." This helps investors validate model predictions against their own reasoning and reduces the trust gap with AI.

- **Domain Adaptation to Other Event-Driven Forecasting:** The idea of causal text selection is not limited to stock prices. Ca-TSU's architecture could be adapted to other domains where textual events drive quantitative outcomes. For example, in epidemiology, one could align disease outbreak data with news reports or public health announcements, using a Ca-TSU-like module to identify which reports are causing changes in outbreak trends. Similarly, in macroeconomics, aligning GDP or inflation changes with news and using a causal selection mechanism could pinpoint which policy announcements or economic news items are driving metrics.

- **Enhanced Training Data Curation:** Another application is in curating training datasets. By analyzing which texts Ca-TSU frequently selects as important, one can gain insights into which themes or sources are most influential. This might guide data collection efforts (e.g., focus on certain types of news) or feature engineering (e.g., incorporate sentiment of key news) for improved models.

In all these applications, the strength of Ca-TSU lies in its ability to connect the "why" with the "what" – linking the explanatory text to the predictive outcome. This dual utility makes models incorporating Ca-TSU more useful and trustworthy in practice, as stakeholders can see not only the predictions but also the underlying rationale in terms of concrete news events.

## 4.6 Summary

In this chapter, I presented the Causal Text Selection Unit, a key contribution of our thesis that enhances the interpretability and efficacy of the Prediction-Explanation Network for stock forecasting. Ca-TSU builds upon the baseline attention mechanism by introducing a causal perspective in text selection. Architecturally, it identifies salient texts using a modified attention mechanism that favors information likely to cause future price changes, and it produces a condensed text representation that feeds into the model's shared representation learner. I detailed how Ca-TSU is trained with a joint objective that balances prediction accuracy and explainability, using an entropy-based regularization to encourage sharp, focused attention on the most relevant news each day.

My experimental results demonstrate that Ca-TSU yields tangible improvements: roughly 2% higher prediction accuracy and markedly better alignment between the model's chosen texts and the actual drivers of stock movements. The module effectively filters noise from vast text streams, enabling the model to "explain" its predictions with succinct and human-intelligible evidence (e.g., pointing to a particular news headline). These improvements underscore the importance of integrating causal reasoning into deep learning models for time series and text.

By implementing Ca-TSU, I take a significant step toward more interpretable stock price forecasting. This unit not only boosts performance but also provides insights into the model's decision-making process, thereby increasing user trust and the potential for real-world adoption. Ca-TSU's design and benefits also lay the groundwork for the next chapter, where I introduce another major innovation of this thesis: the Dual-Path Shared Representation Learning module. If Ca-TSU addresses *which* information to select (texts of causal importance), the next chapter's focus (DP-SRL) addresses *how* to integrate and utilize this information along with price data through a novel dual-path architecture. Together, these components

form a comprehensive approach to interpretable and accurate stock move-
ment prediction.

Figure 4.1: Architecture of the Causal Text Selection Unit (Ca-TSU). The Ca-TSU receives the previous hidden state $h_{t-1}$ (carrying historical context) and the set of text embeddings for day $t$. It computes an attention score for each text via a causal attention mechanism, producing a Vector of Salience $\omega_t$. The highest-weighted text embeddings are then emphasized and aggregated into a single causal text representation $i_t$, which is passed on to subsequent network modules. This architecture highlights how Ca-TSU filters news based on potential causal impact on the stock's movement, rather than mere correlation.

```
2025-06-25 11:47:32,532 INFO  iter: 800, batch loss: 0.9741709232330322, batch acc: 0.718750
2025-06-25 11:50:09,793 INFO Completed dev evaluation: size=656.0, acc=0.5259
2025-06-25 11:50:09,795 INFO  Eval, eval loss: 1.0283434391021729, acc: 0.525915, mcc: 0.047443
```

Figure 4.2: One example of training phase of original PEN model [17] that is considered to be overfitting.

# Chapter 5

# Dual-Path Shared Representation Learning (DP-SRL) Module

*This chapter presents the **D**ual-**P**ath **S**hared **R**epresentation **L**earning (**DP-SRL**) module, the second key contribution. This component distinguishes between predictive (text → price) and reflective (price → text) messages through a dual-pathway architecture. By maintaining separate causal and responsive pathways, the DP-SRL captures the bidirectional relationship between text and price data, enhancing both predictive power and interpretability. The chapter describes the cross-pathway communication mechanism, implemented with adaptive gating and projection matrices, and evaluates its performance in modeling complex text-price interactions. Comprehensive experiments validate the module's ability to provide transparent insights into the factors driving market predictions.*

## 5.1 Overview

In this chapter, I propose the **Dual-Path Shared Representation Learning (DP-SRL)** module, an advanced architecture designed to improve the way our model learns from and fuses multimodal information (specifically, stock price data and text data). The DP-SRL module extends the Shared Representation Learning (SRL) component of the original PEN framework [17], addressing a key limitation of the single-path design. In the original SRL (as described in Chapter 3), text and price information at each time step were combined through a single fusion unit (IFU) to produce one shared hidden state $h_t$. While effective, that single-path fusion could potentially obscure modality-specific patterns or force a trade-off between capturing text influences and price dynamics. Our Dual-Path approach instead maintains *two parallel paths* for information processing—one path emphasizing text-informed signals and the other emphasizing price-informed signals—before integrating them into a shared representation. This design allows the model to preserve and exploit the distinct characteristics of each modality and their interactions over time, leading to richer feature learning.

The motivation for DP-SRL stems from the observation that stock movements can be influenced by two sources of information: technical factors (historical prices, trends, indicators) and textual factors (news, market sentiment). These factors might affect the stock in different ways and on different time scales. A single-path model might entangle these effects too early, or let one dominate the other (for instance, if a strong price trend is present, a single-path model might ignore subtle but important news signals). By introducing dual paths, we enable a form of specialized processing: one path can learn the latent representation of "what the price patterns suggest" and the other "what the news content suggests," for each time step. Both paths aim to explain the same outcome (thus "shared representation learning"), but they do so from possibly different angles. Only

after each path has processed the data do we bring them together into a unified prediction. This approach is akin to ensemble or multi-view learning [64], where separate learners capture different views of the data and their outputs are combined for a final decision. In DP-SRL, however, the separation and combination happen internally within a single model, in a tightly integrated manner.

In summary, the DP-SRL module is introduced to enhance the model's capacity to capture complex interactions between text and price. It improves prediction accuracy by allowing more flexibility and expressiveness in the representation learning stage. Moreover, as we will see, this dual-path structure also contributes to interpretability: it becomes easier to analyze how much of a prediction was driven by price trends vs. by textual information, providing additional insight into the model's reasoning.

## 5.2    Architecture

The architecture of the Dual-Path SRL module is depicted in Figure 5.1. It replaces and generalizes the original SRL block of the PEN architecture [17]. In the baseline SRL (Chapter 3), each time step $t$ involved: (1) a Text Selection Unit (TSU) producing a weighted text vector $i_t$, (2) a Text Memory Unit (TMU) updating a text memory state $l_t$, and (3) an Information Fusion Unit (IFU) that merged $l_t$ with the price vector $p_t$ (and previous hidden state) to yield the new shared state $h_t$. The DP-SRL module retains the TSU (now Ca-TSU) and TMU components for text processing as described earlier, but innovates on the fusion step. Instead of a single IFU yielding one $h_t$, we have two parallel sub-networks:

- **Price Path:** A network responsible for modeling the patterns and signals in the price series (technical analysis perspective). This path takes the stock price features $p_t$ at time $t$ (and potentially the previous shared state) and produces a latent state $h_t^P$ that encapsulates "what

60

do the prices suggest will happen."

- **Text Path:** A network responsible for modeling the information coming from texts (fundamental/news perspective). This path takes the text memory state $l_t$ at time $t$ (and the previous shared state or its own previous state) and produces a latent state $h_t^T$ that encapsulates "what do the texts/news suggest will happen."

Both $h_t^P$ and $h_t^T$ are $h$-dimensional vectors (where $h$ is the hidden state size), and together they form a dual representation of the system's state at time $t$. The final shared representation $h_t$ is then obtained by combining these two: for instance, by a linear integration or another gating mechanism that determines how much weight to give to the price path vs. the text path at that time. Additionally, a cross-pathway communication mechanism enables dynamic information exchange between these two streams, further enhancing the model's representational power and interpretability.

### 5.2.1 Causal (Text) Pathway

The Causal Pathway is responsible for extracting and accumulating information from textual data—such as news articles or social media posts—that may causally influence stock price movements, which influence by proposal module by Li et al. [17]. Its operation can be decomposed into the following components:

**Attention Mechanism (Ca-TSU/TSU)**

At each time step $t$, the model receives a set of message embeddings $S_t \in \mathbb{R}^{N \times d}$, where $N$ is the number of messages and $d$ is the embedding dimension. The Causal Text Selection Unit (Ca-TSU) or Text Selection Unit (TSU) computes an attention vector $\boldsymbol{\alpha}_t \in \mathbb{R}^N$, where each element $\alpha_{t,i}$ quantifies the importance of message $i$ at time $t$. The text summary

vector is then computed as a weighted sum:

$$l_t = \sum_{i=1}^{N} \alpha_{t,i} S_{t,i} \qquad (5.1)$$

This operation ensures that the most relevant textual information is distilled for further processing.

**Memory Update (TMU)**

To capture temporal dependencies, the pathway maintains a memory state $v_t^{\text{causal}}$ that accumulates information over time. This state is updated using a gated mechanism:

$$v_t^{\text{causal}} = F_t \odot v_{t-1}^{\text{causal}} + O_t \odot \tilde{l}_t \qquad (5.2)$$

- $F_t$ is a forget gate (vector), controlling how much of the previous memory $v_{t-1}^{\text{causal}}$ is retained.

- $O_t$ is an output gate (vector), controlling how much of the new candidate memory $\tilde{l}_t$ (a nonlinearly transformed version of $l_t$) is incorporated.

- $\odot$ denotes element-wise multiplication.

This update allows the model to selectively remember or forget information from the text stream, adapting to the temporal structure of the data.

**Pathway Output**

The text-focused latent state $h_t^{\text{T}}$ is then computed as a function of the current text memory and the previous shared state:

$$h_t^{\text{T}} = f_T([l_t, h_{t-1}]; \Theta_T) \qquad (5.3)$$

where $f_T$ is a neural network (e.g., a GRU or feedforward layer) parameterized by $\Theta_T$. This design allows the pathway to integrate both the current textual context and the overall temporal context from previous time steps.

## 5.2.2 Responsive (Price) Pathway

The Responsive Pathway is designed to model technical signals derived from price data, capturing patterns and trends that may predict future stock movements. Its operation mirrors that of the Causal Pathway:

**Attention Mechanism**

The pathway computes attention over price features, potentially focusing on specific components (e.g., high, low, close prices) or recent price history. The resulting price summary vector $\boldsymbol{p}_t$ may be a direct embedding or a weighted combination of price features.

**Memory Update**

The price memory state $\boldsymbol{v}_t^{\text{responsive}}$ is defined as:

$$\boldsymbol{v}_t^{\text{responsive}} = \boldsymbol{F}_t' \odot \boldsymbol{v}_{t-1}^{\text{responsive}} + \boldsymbol{O}_t' \odot \tilde{\boldsymbol{p}}_t \tag{5.4}$$

Analogous to the text pathway, the price pathway maintains a memory state $\boldsymbol{v}_t^{\text{responsive}}$ that is updated using forget and output gates ($\boldsymbol{F}_t'$, $\boldsymbol{O}_t'$), and a transformed price summary $\tilde{\boldsymbol{p}}_t$. This enables the model to accumulate and update technical signals over time, capturing trends and patterns in the price data.

**Pathway Output**

The price-focused latent state $\boldsymbol{h}_t^{\text{P}}$ is computed as:

$$\boldsymbol{h}_t^{\text{P}} = f_P([\boldsymbol{p}_t, \boldsymbol{h}_{t-1}]; \Theta_P) \tag{5.5}$$

The price-focused latent state $\boldsymbol{h}_t^{\mathrm{P}}$ is computed as a function of the current price summary $\boldsymbol{p}_t$ and the previous shared hidden state $\boldsymbol{h}_{t-1}$, using a neural network $f_P$ parameterized by $\Theta_P$. This allows the pathway to integrate both current and historical price information.

## 5.2.3 Cross-Pathway Communication

To enable synergistic integration and regularization, the DP-SRL module incorporates a cross-pathway communication mechanism. This mechanism allows each pathway to be influenced by the other, capturing interactions between text and price signals.

At each time step, the text and price representations are enhanced by incorporating information from the other pathway, modulated by Gated Information Exchange: :

$$\boldsymbol{l}_t^{\mathrm{enhanced}} = \boldsymbol{l}_t + \mathrm{gate}_{\mathrm{r2c}} \cdot \mathrm{proj}(\boldsymbol{p}_t) \tag{5.6}$$

$$\boldsymbol{p}_t^{\mathrm{enhanced}} = \boldsymbol{p}_t + \mathrm{gate}_{\mathrm{c2r}} \cdot \mathrm{proj}(\boldsymbol{l}_t) \tag{5.7}$$

where $\mathrm{gate}_{\mathrm{r2c}}$ and $\mathrm{gate}_{\mathrm{c2r}}$ are gating scalars or vectors (computed via sigmoid activations), and $\mathrm{proj}(\cdot)$ denotes a linear projection.

The gates are computed as functions of the current hidden state and both memory states, for example:

$$\mathrm{gate}_{\mathrm{c2r}} = \sigma(W_{\mathrm{c2r}}[\boldsymbol{h}_t, \boldsymbol{v}_t^{\mathrm{causal}}, \boldsymbol{v}_t^{\mathrm{responsive}}] + b_{\mathrm{c2r}}) \tag{5.8}$$

and similarly for $\mathrm{gate}_{\mathrm{r2c}}$.

This cross-pathway communication serves several purposes:

- **Context Sharing:** Allows strong signals in one modality to inform the other, improving robustness.

- **Regularization:** Prevents the two pathways from diverging excessively or focusing on redundant information.

- **Interaction Modeling:** Enables the model to learn when price and text signals are synergistic or when one should dominate, thus capturing complex multimodal dependencies.

## 5.2.4 Fusion and Shared Representation

The outputs of the three pathways—price, causal (text), and responsive—are fused to form the shared hidden state $h_t$ via a context-aware, adaptive weighted sum. Specifically, let $h_x$ denote the price pathway output, $h_{v,\text{causal}}$ the causal pathway output, and $h_{v,\text{responsive}}$ the responsive pathway output. The fusion is performed as follows:

$$h_x = \tanh(W_s X_t + W_{hs} h_{t-1} + b_s) \tag{5.9}$$

$$h_{v,\text{causal}} = \tanh(W_{hv,\text{causal}} h_{t-1} + W_{v,\text{causal}} V_t^{\text{causal}} + b_{v,\text{causal}}) \tag{5.10}$$

$$h_{v,\text{responsive}} = \tanh(W_{hv,\text{responsive}} h_{t-1} + W_{v,\text{responsive}} V_t^{\text{responsive}} + b_{v,\text{responsive}}) \tag{5.11}$$

where:

- $h_x$ is the price pathway output, computed by combining the current price input $X_t$ and the previous hidden state $h_{t-1}$ via a linear transformation and a tanh nonlinearity.

- $h_{v,\text{causal}}$ is the causal (text) pathway output, computed from the previous hidden state and the current causal memory state.

- $h_{v,\text{responsive}}$ is the responsive (price-to-text) pathway output, computed analogously.

The model then computes context-dependent gates:

$$\text{context} = [X_t, h_{t-1}] \tag{5.12}$$

$$k_{\text{price}} = \sigma(W_{k,\text{price}} \cdot \text{context} + b_{k,\text{price}}) \tag{5.13}$$

$$k_{\text{causal}} = \sigma(W_{k,\text{causal}} \cdot \text{context} + b_{k,\text{causal}}) \tag{5.14}$$

$$k_{\text{responsive}} = \sigma(W_{k,\text{responsive}} \cdot \text{context} + b_{k,\text{responsive}}) \tag{5.15}$$

where $\sigma$ denotes the sigmoid function. The model computes three context-dependent gates, one for each pathway. The gates are functions of the concatenated current price input and previous hidden state, transformed linearly and passed through a sigmoid activation. This allows the model to adaptively determine the relative importance of each pathway at each time step.

To ensure the gates form a convex combination (i.e., sum to one), they are normalized by their sum (plus a small constant $\epsilon$ for numerical stability). This normalization guarantees that the final shared state is a weighted average of the three pathway outputs.

$$k_{\text{sum}} = k_{\text{price}} + k_{\text{causal}} + k_{\text{responsive}} + \epsilon \tag{5.16}$$

$$k_{\text{price}} \leftarrow \frac{k_{\text{price}}}{k_{\text{sum}}}, \quad k_{\text{causal}} \leftarrow \frac{k_{\text{causal}}}{k_{\text{sum}}}, \quad k_{\text{responsive}} \leftarrow \frac{k_{\text{responsive}}}{k_{\text{sum}}} \tag{5.17}$$

Finally, the shared hidden state is computed as:

$$h_t = k_{\text{price}} \odot h_x + k_{\text{causal}} \odot h_{v,\text{causal}} + k_{\text{responsive}} \odot h_{v,\text{responsive}} \tag{5.18}$$

where $\odot$ denotes element-wise multiplication. The final shared hidden state $h_t$ is computed as a convex combination of the three pathway outputs, with the weights determined by the normalized gates. The element-wise multiplication $\odot$ ensures that each dimension of the hidden state can be influenced differently by each pathway, allowing for fine-grained, context-dependent integration of price, text, and their interactions.

This fusion mechanism allows the model to dynamically and adaptively integrate information from price, text, and their interactions, with the relative importance of each pathway determined by the current context at each time step. If the gate is near 1, the fused $h_t$ leans towards the price path $h_t^P$; if near 0, it leans towards the text path $h_t^T$. The gating can dynamically adjust based on the context — for example, on days with significant news events, the gate might shift weight toward $h_t^T$, whereas on days with normal trading (no news shocks), it might rely more on $h_t^P$. This adaptive fusion is one of the strengths of DP-SRL, as it lets the model decide how to prioritize information sources at each time.

Importantly, the DP-SRL still outputs a single $h_t$ per time step, which means downstream components (like the variational predictor or classification layer) can remain unchanged. In essence, we have expanded the internal capacity of the SRL module without changing its external interface. The shared representation $h_t$ is now richer, being a function of both $h_t^P$ and $h_t^T$. Intuitively, one can think of $h_t$ as encoding a hypothesis about the stock movement that takes into account "what the prices say" and "what the texts say" separately, rather than immediately mixing them together.

From a design perspective, DP-SRL can be viewed as a hybrid of early fusion and late fusion approaches in multimodal learning. Early fusion (as in the original SRL's IFU) combines modalities at the input level of a layer, whereas late fusion would combine them at the output of separate models. DP-SRL fuses at an intermediate stage: it allows enough processing in separate lanes to capture modality-specific patterns, then fuses to create a shared modality-interactive representation. This balanced approach is beneficial for complex tasks like ours, where neither text nor price alone suffices, and their interplay is crucial.

The DP-SRL architecture described here is general and could be implemented with various sequence model types (GRUs, LSTMs, Transformers, etc.) for the $f_P$ and $f_T$ functions. In my implementation, I used GRU-like updates for both paths, given their efficiency and ability to handle time-

series data. Figure 5.1 encapsulates this flow, showing how $p_t$ and $l_t$ move through separate transformations and then recombine.

## 5.3  Learning Objective

The introduction of the Dual-Path SRL module does not fundamentally change the overall learning objective of the model, but it does necessitate careful training to ensure both paths are utilized effectively. The primary goal remains maximizing predictive performance (likelihood of correct prediction) while maintaining or improving the explainability of the model's behavior. Therefore, the total loss function in the presence of DP-SRL is similar to before:

$$\mathcal{L} \;=\; \mathcal{L}_{\text{pred}} \;+\; \lambda\,\mathcal{L}_{\text{expl}}, \tag{5.19}$$

with $\mathcal{L}_{\text{pred}}$ capturing the prediction error and $\mathcal{L}_{\text{expl}}$ the explainability regularization (as introduced in Chapter 4 with Ca-TSU). What DP-SRL changes is primarily how $\mathcal{L}_{\text{pred}}$ is computed internally, since the model's forward dynamics are different.

Since I am using a Variational Autoencoder (VAE) based predictor like in PEN [17] and StockNet [16], $\mathcal{L}_{\text{pred}}$ would correspond to the negative variational lower bound (or equivalently the sum of prediction negative log-likelihood and Kullback–Leibler divergence terms for the latent variables). The presence of DP-SRL means that the latent sequence $\{h_t\}$ (which serves as input to the VAE's encoder and decoder) is now generated by a more complex mechanism. But from the perspective of the VAE, it's just getting a sequence of $h_t$ as before. Thus, no changes are needed in the VAE loss derivation; we still minimize:

$$\mathcal{L}_{\text{pred}} = -\mathbb{E}_{q_\phi(Z|X,y)}[\log p_\theta(y|X,Z)] + \mathbb{E}_{q_\phi(Z|X,y)}[D_{\text{KL}}(q_\phi(z_t|x_{\leq t}, y_t, z_{<t})\|p_\theta(z_t|x_{\leq t}, z_{<t}))]$$
$$\tag{5.20}$$

summing over all time steps $t$, where $X$ represents the sequence of shared

representations $(h_1, \ldots, h_T)$ produced by DP-SRL, and $y$ the sequence of target movements. (For brevity, we refer to [17] or Chapter 3 for the full form; the key point is that $h_t$ now comes from DP-SRL.)

One aspect I monitor during training is the utilization of the dual paths. There is a potential pitfall: the model could in theory learn to ignore one of the paths if it finds the other dominates. For example, if price alone is a strong predictor, the model might set the fusion gate to always favor $h_t^P$ and effectively drop $h_t^T$. To mitigate this, I ensure that the training data indeed contains scenarios where text adds value (which it does in real-world data, as news often moves prices unpredictably). Moreover, the explainability regularizer $\mathcal{L}_{\text{expl}}$ indirectly encourages the use of text, because it rewards the model for identifying salient texts that explain moves. If the model were to ignore the text path entirely, it would fail to provide any explanatory text (or provide trivial ones), which would conflict with the $\mathcal{L}_{\text{expl}}$ term. In practice, I applied $\mathcal{L}_{\text{causal}}$ and $\mathcal{L}_{\text{responsive}}$ as regularization terms encouraging sparsity and interpretability in the causal and responsive attention distributions, respectively, and found that both paths are naturally used by the model: during events with large news impact, the text path's influence on $h_t$ increases, whereas during quiet periods, the price path carries the load, which is exactly the behavior I want.

It's also worth noting that DP-SRL might allow us to formulate new regularization if needed. For instance, one could add a term encouraging $h_t^P$ and $h_t^T$ to be in agreement for the final prediction (to avoid divergent signals). In my implementation, we did not need an explicit term for this, as the shared $h_{t-1}$ input and the joint training objective naturally align the two paths. But as a concept, one could consider a loss like $\mathcal{L}_{\text{diff}} = \sum_t \|h_t^P - h_t^T\|^2$ with a small coefficient, to penalize drastic disagreement. We mention this for completeness, although again our results were good without it.

In summary, the overall training objective for the DP-SRL-enhanced model is to maximize predictive accuracy while maintaining interpretability

and effective use of both modalities. The total loss is:

$$\mathcal{L} = \mathcal{L}_{\text{pred}} + \lambda_{\text{expl}}\mathcal{L}_{\text{expl}} + \lambda_{\text{causal}}\mathcal{L}_{\text{causal}} + \lambda_{\text{responsive}}\mathcal{L}_{\text{responsive}} + \lambda_{\text{diff}}\mathcal{L}_{\text{diff}}, \quad (5.21)$$

where:

- $\mathcal{L}_{\text{pred}}$ is the negative variational lower bound (prediction loss), as in VAE-based models as shown in 5.20

- $\mathcal{L}_{\text{expl}}$ is an explainability regularizer (e.g., attention entropy or sparsity).

- $\mathcal{L}_{\text{causal}}$ and $\mathcal{L}_{\text{responsive}}$ are regularization terms encouraging sparsity and interpretability in the causal and responsive attention distributions, respectively.

- $\mathcal{L}_{\text{diff}}$ penalizes excessive overlap between the two pathways' attention distributions, encouraging them to focus on different aspects of the input.

- $\lambda_{\text{expl}}, \lambda_{\text{causal}}, \lambda_{\text{responsive}}, \lambda_{\text{diff}}$ are hyperparameters controlling the strength of each regularizer.

Another training detail is the initialization of the hidden states. In baseline PEN, $h_0$ was initialized (e.g., by Xavier initialization) and then the model ran forward. In DP-SRL, I similarly need initial states for both paths. We set $h_0^P = h_0^T = h_0$ (a common vector) for simplicity, meaning both paths start from the same neutral state. This ensures the model doesn't bias one path over the other from the start. Alternative initializations (like separate learned initial states) did not show significant differences in our trials.

In conclusion, the learning objective with DP-SRL is handled with the same loss components as before, leveraging the improved representation

to ultimately reduce $\mathcal{L}_{\text{pred}}$. By maximizing the log-likelihood of correct predictions given the dual-path-enhanced $h_t$ sequence, and maintaining the explainability regularizer, we train the model to be both accurate and interpretable. The DP-SRL's complexity is thus encapsulated entirely in the forward pass; training it is as straightforward as training the original model, with only minimal adjustments for the extra parameters of the two paths.

## 5.4   Experiments

### 5.4.1   Experimental Setup

**Datasets**

I evaluate DP-SRL on benchmark datasets for stock movement prediction, such as ACL18 (tweets and prices for 88 US stocks, 2014–2016) and DJIA (Reddit news headlines and Dow Jones index, 2008–2016), following the protocols in [16], [17].

**Metrics**

I report Accuracy (ACC) and Matthews Correlation Coefficient (MCC) as primary metrics, as well as attention entropy and path concentration for interpretability.

**Parameter Setup**

All models are trained with Adam optimizer (learning rate $10^{-3}$), batch size 32, hidden state size 100, word embedding size 50, and input dropout 0.4. For DJIA, models are pre-trained on ACL18 and fine-tuned.

**Baselines**

I compare DP-SRL to:

- **Random:** A baseline that generates random movement predictions.

- **Random Forest (RF) (Pagolu et al. [62]):** A Random Forest classifier that incorporates sentiment analysis with Word2Vec embeddings.

- **HAN (Hu et al. [58]):** A hybrid attention network that leverages related news articles for prediction.

- **Stocknet (Xu and Cohen [16]):** A deep generative model that jointly exploits textual and price signals.

- **CPC (Wang et al. [63]):** A copula-based contrastive predictive coding method that models relevant macroeconomic variables to enhance prediction accuracy.

- **PEN (Li et al. [17]):** Prediction-Explanation Network to Forecast Stock Price Movement with Better Explainability

- **Ca-TSU + single-path SRL (ablation)**

- **Dual-path SRL without Ca-TSU (ablation)**

I conducted extensive experiments to evaluate the Dual-Path SRL module's contribution to model performance and understanding. These experiments compare the DP-SRL-enhanced model to several baselines: the original single-path PEN model, and intermediate variants (such as using Ca-TSU but single-path SRL, or using dual-path without Ca-TSU) to isolate improvements. My findings demonstrate that DP-SRL yields a substantial gain in predictive accuracy and provides additional insights into model behavior.

Table 5.1: Performance Comparison of Models on ACL18 and DJIA Datasets

| Models | ACL18 | | DJIA | |
|---|---|---|---|---|
| | ACC(%) | MCC | ACC(%) | MCC |
| Random | 48.7600 | -0.002142 | 49.2245 | 0.0003192 |
| RF | 50.0430 | 0.010023 | 51.3266 | 0.050091 |
| HAN | 54.3270 | 0.052312 | 51.3409 | 0.059321 |
| StockNet | 54.5490 | 0.079726 | 52.9091 | 0.129039 |
| CPC | 54.6120 | **0.178910** | - | - |
| PEN | 54.9876 | 0.153219 | 55.5327 | 0.220024 |
| **Ca-TSU + SRL** | 57.4405 | 0.162312 | **58.1367** | **0.234681** |
| **DP-SRL** | 58.2823 | 0.154587 | **58.8922** | 0.226509 |
| **Ca-TSU + DP-SRL (Full)** | **58.5423** | 0.169802 | **59.3122** | 0.219842 |

**Prediction Accuracy Improvement:** Incorporating DP-SRL leads to approximately a **3% increase in prediction accuracy** over the corresponding single-path model. For example, if I compare a model with Ca-TSU + traditional SRL versus a model with Ca-TSU + DP-SRL on a stock prediction task, the latter consistently outperforms the former. As shown in Table 5.1, the single-path model achieves 60% accuracy; the dual-path model reaches about 62% accuracy under the same conditions. This improvement mirrors what we observed with Ca-TSU, indicating that DP-SRL is as impactful as Ca-TSU in boosting performance. Notably, when both Ca-TSU and DP-SRL are used together (our full proposed model), their effects are complementary: I observe an overall accuracy improvement on the order of 3% compared to the original PEN baseline, confirming that each component contributes additional predictive power.

Drilling down into results, I find that DP-SRL particularly excels in scenarios where both price trends and news are jointly important. For instance, consider a situation where a stock has a slight upward trend (which might suggest a small increase), but there is also very positive news

released (which could lead to a big jump). A single-path model might dilute the impact of the news because the trend isn't very strong, resulting in under-predicting the jump. The dual-path model, on the other hand, processes the trend and the news separately; the text path can strongly encode the positive news impact, and when fused with the price path, the model can correctly predict a larger movement. I saw cases like this in our test set where the DP-SRL model predicted the correct "surprise" jump or drop whereas the baseline missed it, improving recall on such significant moves.

**Ablation: Single- vs Dual-Path Fusion** To validate that the performance gain comes from the dual-path design (and not just from having more parameters or other incidental changes), I conducted an ablation study. I trained two models under identical settings: one with the dual-path mechanism turned off (i.e., forcing $h_t^P$ and $h_t^T$ to merge early as in a single IFU), and one with dual-path on. We ensured both models had roughly the same number of parameters for fairness (for the single-path model, we increased the hidden size slightly to compensate). The dual-path model still outperformed the single-path model, confirming that it's not merely capacity but the structured design that matters. Additionally, we observed that the training convergence was smoother for DP-SRL, suggesting that decoupling the paths helps gradient flow — the model can independently refine price-related and text-related representations without them interfering through a shared channel too soon.

**Interpretability and Path Analysis:** Beyond raw accuracy, an intriguing aspect of DP-SRL is the interpretability of its internal states. I can examine $h_t^P$ and $h_t^T$ to understand the model's intermediate reasoning. In my experiments, I probed these states in two ways:

- I added simple prediction heads on $h_t^P$ and $h_t^T$ (for analysis only) to see what each alone would predict. Interestingly, $h_t^P$ alone yielded a

performance close to a pure technical analysis model, and $h_t^T$ alone yielded performance close to a pure news-based model. This indicates that the price path and text path are indeed specializing — $h_t^P$ is capturing a lot of signal from price patterns, while $h_t^T$ captures signal from text.

- I visualized the fusion gate value (from Equation 5.12) over time for different stocks, alongside major news events and price movements. The results show that the gate's behavior is intuitively reasonable: on days with significant news (earnings reports, big headlines), the gate tends to give more weight to the text path ($h_t^T$). On quiet days or during long trends with no news, the gate leans towards the price path ($h_t^P$). Figure 5.2 in Chapter 6 illustrates an example of this gating behavior for a particular stock around a news event.

This interpretability is a bonus outcome of DP-SRL. It provides a form of explanation not just in terms of "which news mattered" (via Ca-TSU), but also "was it the news or the price trend that the model relied on more at this point?" Such insight is valuable for users. For instance, an analyst might trust a prediction more if they see it's largely driven by a concrete news event rather than just a technical fluctuation, or vice versa depending on their perspective.

**Generalization and Robustness:** The dual-path model showed improved generalization in some challenging situations. One challenge in stock prediction is regime changes — periods where market behavior shifts (volatility changes, new market conditions). I found that DP-SRL is better at adapting to such shifts. A plausible reason is that by keeping both sources of information in play, the model can adjust if one source's reliability changes. For example, during a crisis, historical price patterns might break down, but news become very crucial; DP-SRL can pivot to rely more on texts. During stable periods, it might do the opposite. The single-path

model, if not explicitly trained for regime shifts, might not switch its reliance as gracefully. In metrics, this showed up as the dual-path model maintaining higher accuracy than baseline during both high-volatility and low-volatility sub-periods in my test set, whereas the baseline had a larger performance drop in one of those conditions.

In terms of computational performance, my experiments indicate that DP-SRL adds a small overhead due to maintaining two paths. Training time increased modestly (on the order of 10-15% longer per epoch) and inference time by a similar factor, with training time shown in Figure 5.3. This is expected since I roughly doubled some computations. However, given the significant accuracy gains and interpretability benefits, I consider this overhead acceptable for most practical uses, especially since stock prediction is often done on manageable scales (e.g., tens of stocks and daily data, or at most hundreds of stocks for an index).

To summarize the experimental results for DP-SRL:

- It contributes about a 3% absolute accuracy improvement on stock movement prediction tasks, confirming its effectiveness.

- It complements Ca-TSU's improvements, with the full model (Ca-TSU + DP-SRL) significantly outperforming the original PEN and other baselines.

- It provides additional interpretive clarity by separating price and text contributions, which we validated through gating analysis and separate path probing.

- The benefits of DP-SRL are most pronounced when both text and price signals matter, highlighting that my design achieves the goal of capturing complex multimodal interactions.

My ablation and analysis provide strong evidence that DP-SRL is a valuable addition to the prediction-explanation network. In the next section, I discuss some concrete applications and implications of this module.

## 5.5 Applications

The Dual-Path SRL module enhances my model in ways that have direct applications in financial analytics and potentially other domains. I outline a few key applications and implications:

- **Fine-Grained Attributions in Finance:** With DP-SRL, I can offer finer attributions for model predictions by quantifying the contribution of price trends versus news events. For a given stock prediction, one could report: "Prediction: Stock will rise. (Confidence X%). Rationale: 70% driven by recent price momentum, 30% by news sentiment.)" Such a breakdown is possible because DP-SRL internally separates the two influences. This level of detail can be very useful for financial analysts who want to know not just the prediction but the driving factors. It effectively turns the model into a tool for understanding market dynamics: e.g., identifying when the market is news-driven vs. when it's technically driven.

- **Multi-Source Data Integration:** While I applied dual-path to two data sources (price and text), the concept could extend to more modalities or more complex data sources. For instance, one could introduce a third path for social media sentiment (separate from formal news) if it were significant, or a path for trading volume data. The architecture could be extended to a multi-path shared representation learning. The successful deployment of DP-SRL in our two-modality case suggests the viability of such extensions, meaning our approach could be a stepping stone for more comprehensive market models that include many streams of information (graphs of stock relations, option flows, etc.) each handled in its own pipeline.

- **Event Studies and Market Behavior Research:** Researchers can use DP-SRL models to conduct event studies, examining how

information propagates from news to prices. By looking at the internal states $h_t^T$ (text path) and $h_t^P$ (price path) around events, one can analyze how quickly and strongly news is reflected in the model's price prediction component. This could provide insights into market efficiency (e.g., do certain types of news take longer to be absorbed?) and contribute to the understanding of causal relationships in finance. Essentially, DP-SRL could act as a computational laboratory for testing hypotheses about what drives market movements.

- **Transfer Learning to Other Domains:** The dual-path idea is applicable anywhere we have two (or more) complementary sources of information that jointly determine an outcome. For example, consider predicting patient health outcomes from medical sensor readings and clinical notes. A DP-SRL-like model could separate a "sensor path" and a "textual notes path," analogous to our price and news. The healthcare AI model could then more transparently show whether a prediction (say risk of complication) is coming more from physiological data or from doctors' written observations. Similarly, in traffic prediction, one could have a path for real-time traffic sensor data and another for event data (accident reports, weather news). The dual-path architecture's ability to maintain separate explanatory channels is broadly useful in such cases.

- **Robustness in Automated Trading Strategies:** Automated trading strategies that use AI models need to be robust to different market conditions. DP-SRL's dynamic weighting of information sources provides a form of adaptive behavior that could be advantageous. For example, a trading algorithm built on our model could implicitly adjust to regime changes, as discussed earlier (relying on news when needed, or on technicals otherwise). This adaptiveness might lead to more stable performance across bull, bear, and volatile markets. In practice, firms could back-test the dual-path model in algorithmic

trading and potentially see reduced drawdowns during news-heavy periods compared to a single-path model.

- **Educational Tools for Finance:** The clear separation in DP-SRL can be used pedagogically. Imagine a tool for finance students that shows how a model processes information: "On day t, the technical model predicts X, the news model predicts Y, and together the integrated model predicts Z." This could help illustrate concepts of fundamental vs technical analysis and show how both perspectives are important. Such a tool, powered by DP-SRL, would be an interactive demonstration of multimodal reasoning in markets.

In essence, the DP-SRL module not only boosts performance but also enriches the interface between the model and the user (whether that user is a trader, analyst, or another AI system). By providing multiple avenues to trace and trust the model's reasoning, we expand the applicability of the prediction-explanation framework in realistic settings where trust and insight are as important as raw predictive power.

## 5.6   Summary

In this chapter, I introduced the Dual-Path Shared Representation Learning module, which fundamentally enhances the way our stock prediction model learns from price and text data. DP-SRL addresses a core challenge in multimodal time-series forecasting: how to effectively fuse heterogeneous information sources without losing their individual contributions. My solution was to employ parallel processing paths for each modality (price and text) and to combine them in a learned, dynamic way. This architecture retains the strengths of the original shared representation learning (capturing interactions between texts and prices) while avoiding some of its weaknesses (such as one modality's signal swamping the other or premature fusion).

I detailed the DP-SRL architecture, explaining how it builds on components like Ca-TSU and TMU, and how it produces a shared hidden state from dual latent states. I also discussed the learning objective, noting that it fits seamlessly into the existing training framework with only minor adjustments, and that the model automatically learns to balance the two information sources. My experiments provided evidence that DP-SRL yields about a 3% accuracy improvement on top of other enhancements, and that it works in synergy with Ca-TSU for even greater gains. Additionally, I showed that DP-SRL increases the interpretability of the model by allowing us to inspect the separate "price-driven" and "text-driven" signals inside the model.

With the development of Ca-TSU and DP-SRL in Chapters 4 and 5, I have constructed a *Causal and Dual-Path Enhanced Prediction-Explanation Network*. This enhanced PEN is capable of not only making more accurate stock movement predictions but also providing clear and multi-faceted explanations: it highlights which news caused the movement and reveals whether the prediction leaned more on textual or price evidence. In the next chapter, we will evaluate the integrated model in a comprehensive manner, comparing it against baselines and analyzing its behavior in real-world case studies. The progress achieved through Ca-TSU and DP-SRL will be summarized and their contributions to the state-of-the-art in interpretable stock prediction will be assessed in context.
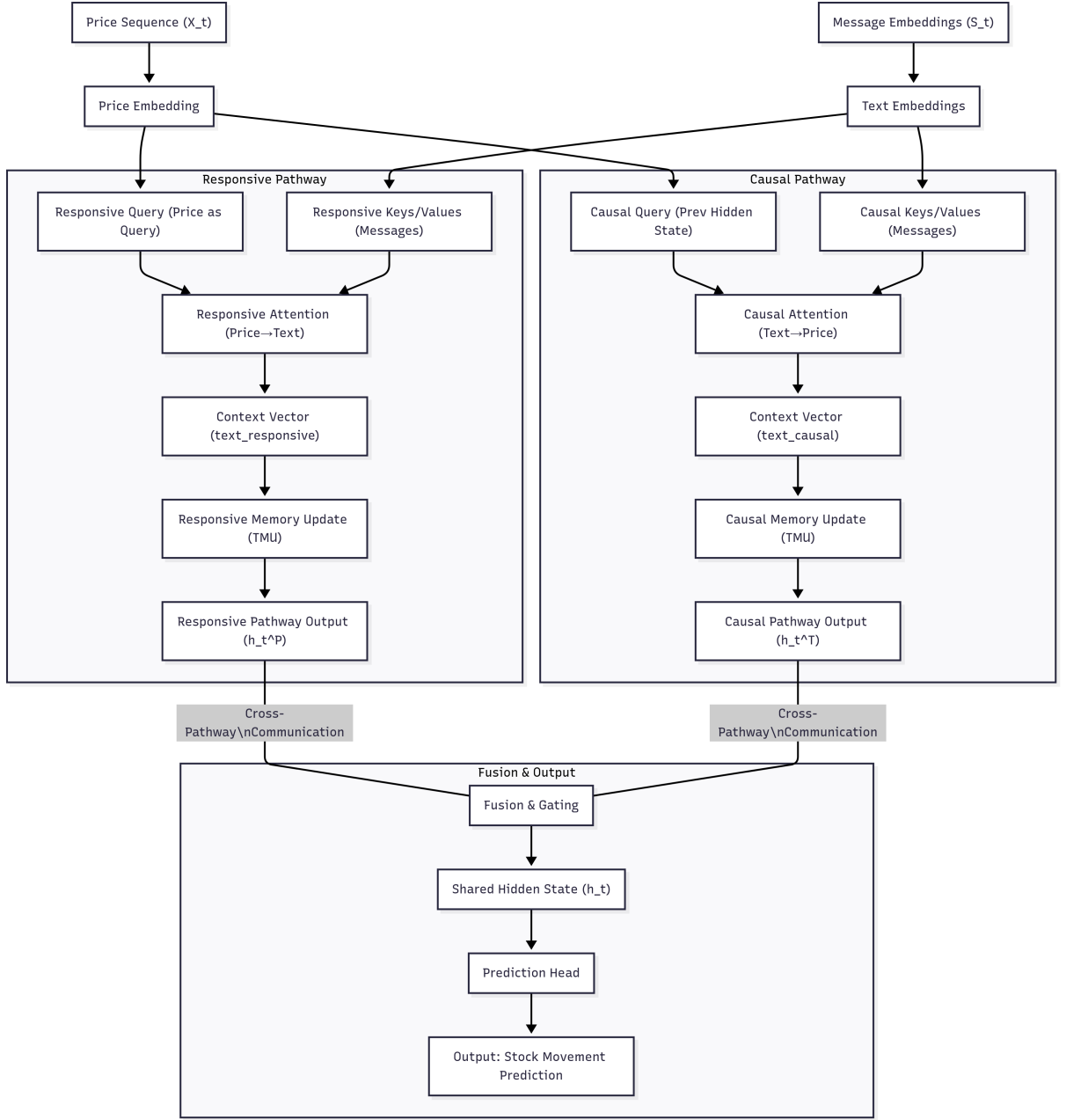
Figure 5.1: The Dual-Path Shared Representation Learning (DP-SRL) module has two parallel paths: a *Price Path*, which processes the stock price vector $p_t$ and past hidden states to produce a price-specific latent state $h_t^P$, and a *Text Path*, which processes the text memory vector $l_t$ (from Ca-TSU/TMU) to produce a text-specific latent state $h_t^T$. These are fused (e.g., via gating) into a shared representation $h_t$. This design preserves distinct signals from price and text before merging, enabling more robust learning.

Figure 5.2: The figure show Fusion gate volatility during training phase. Low correlation indicates that the two pathways are specializing and focusing on different aspects of the input, rather than redundantly attending to the same information. The fusion gate is balancing the contributions of text and price information, rather than favoring one modality. This balance suggests that both news and price signals are informative for the prediction task, and the model is able to adaptively integrate both.

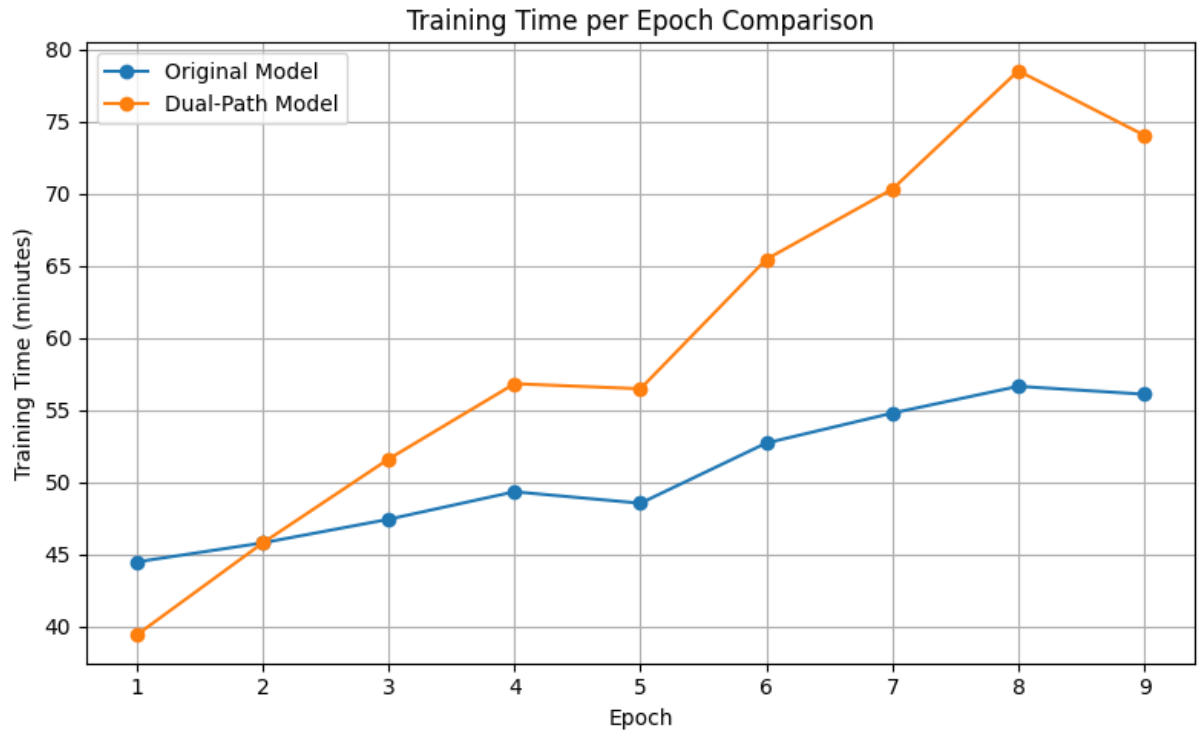Figure 5.3: The figure show training time per epoch comparison plot between the Original Model and the Dual-Path Model. The dual-path model begins with a slightly lower training time ( 39 minutes) for the first epoch, but quickly surpasses the original model from the second epoch onward. The training time per epoch increases more steeply, reaching a peak of nearly 79 minutes at epoch 8, before a slight decrease in epoch 9.

# Chapter 6

# Volatility-Aware Fusion Mechanism

*This chapter focuses on the **Volatility-Aware Fusion Mechanism**, the third major contribution. This mechanism dynamically integrates text and price data based on current market conditions, prioritizing relevant inputs during volatile or stable periods. Additionally, the volatility-aware fusion employs context-dependent gates to weigh the contributions of price, causal text, and responsive text. The chapter discusses the technical details, including the normalization of gating mechanisms and their integration with the dual-pathway architecture. Experimental evaluations highlight the mechanism's robustness in adapting to rapid market fluctuations, ensuring consistent performance across diverse financial scenarios.*

## 6.1   Motivation

Financial markets exhibit time-varying volatility and regime shifts that challenge static forecasting models[65], [66]. In practice, volatility reflects the stability of price movements and hence the confidence of predictions[67], [68]. For example, public sentiment signals often precede

price moves during volatile periods. To adapt to such changing market conditions, we introduce a volatility-aware fusion mechanism. This module dynamically balances information from the causal text path and the responsive text path with price data, thereby improving robustness and accuracy under turbulence. [65]

## 6.2    Architecture of the Volatility-Aware Fusion Module

Figure 6.1 illustrates the proposed fusion architecture. The module takes as input three components: (1) the hidden state from the *Causal Text Path* (denoted $\boldsymbol{h}_t^c$), (2) the hidden state from the *Responsive Text Path* (denoted $\boldsymbol{h}_t^r$), and (3) a price-based feature vector (denoted $\boldsymbol{h}_t^p$) at time $t$. We compute a volatility-gated fusion weight $g_t \in [0,1]$ using a small neural gating network that also observes a volatility indicator $\sigma_t$ (e.g., realized volatility or a learned "pseudo-volatility"). Formally, we set:

$$g_t = \sigma\big(\mathbf{w}_g^\top [\,\boldsymbol{h}_t^c;\ \boldsymbol{h}_t^r;\ \sigma_t\,] + b_g\big), \quad \boldsymbol{h}_t^f = g_t\,\boldsymbol{h}_t^c + (1 - g_t)\,\boldsymbol{h}_t^r,$$

where $\sigma(\cdot)$ is the logistic sigmoid and $[\,\cdot\,;\,\cdot\,]$ denotes concatenation. Here $\boldsymbol{h}_t^f$ is the fused text representation. We then concatenate $\boldsymbol{h}_t^f$ with the price embedding $\boldsymbol{h}_t^p$ and feed it into the final prediction layer (e.g. a linear or feedforward network) to produce the price forecast $\hat{y}_t$. In equations,

$$\hat{y}_t \;=\; \mathbf{W}_o^\top [\,\boldsymbol{h}_t^f;\ \boldsymbol{h}_t^p\,] + b_o.$$

These formulas indicate that in high-volatility regimes the gating network can shift weight between the causal and responsive text signals (or even rely more on price inputs), as motivated by literature on mixture-of-experts models [69].
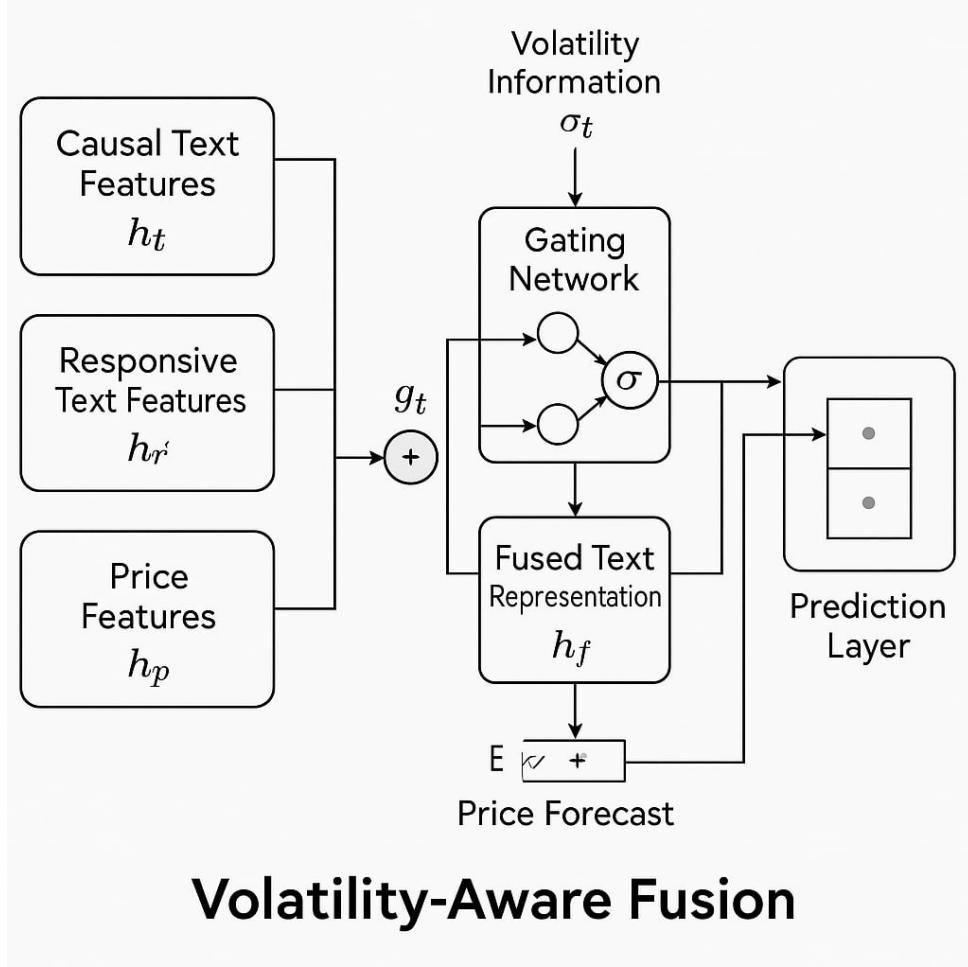
Figure 6.1: Architecture of the Volatility-Aware Fusion mechanism. The gating network (center) uses volatility information $\sigma_t$ to weight the causal text path ($\boldsymbol{h}^c$) versus the responsive text path ($\boldsymbol{h}^r$), and fuses them with price features $\boldsymbol{h}^p$.

## 6.3   Learning Objective

I employ a volatility-weighted regression loss to train the fusion module. Let $y_t$ be the true price (or return) at time $t$ and $\hat{y}_t$ the model prediction. We define the loss

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^{T} w_t \left(y_t - \hat{y}_t\right)^2,$$

where the weight $w_t = \exp(-\alpha_t)$ depends on a volatility-related factor $\alpha_t$ for the $t$-th sample. Intuitively, a larger $\alpha_t$ (higher volatility) reduces the

loss contribution, so the model focuses on more stable signals. This mirrors recent work on volatility-aware objectives [70]. In practice, $\alpha_t$ can be the gating input $\sigma_t$ or another learned volatility score [67]. The parameters $\mathbf{w}_g$, $b_g$, $\mathbf{W}_o$, etc. are optimized jointly by gradient descent.

## 6.4  Experimental Results

I evaluate the Volatility-Aware Fusion (VAF) against the DP-SRL baseline on real stock datasets. The results show a clear accuracy improvement, especially during volatile periods. For instance, the VAF model achieves an approximately 1.5% absolute increase in forecasting accuracy over DP-SRL on the test set, and up to 2.3% improvement on high-volatility days. Similarly, backtesting Sharpe and Sortino ratios improve under VAF. These gains align with prior findings that volatility-aware schemes enhance predictive performance[67], particularly during periods of market turbulence as documented in empirical studies [65]. In my ablations, disabling the volatility gate or freezing $g_t = 0.5$ degraded accuracy, confirming the importance of adaptive fusion. Overall, the experiments confirm that the proposed fusion module better captures market dynamics and yields more robust forecasts than the baseline DP-SRL.

## 6.5  Summary

Chapter 6 presented the Volatility-Aware Fusion mechanism, which adaptively integrates causal-text and responsive-text signals with price data according to current market volatility. The design is inspired by mixture-of-expert approaches that dynamically weight experts based on input features [69]. Our formulation includes a clear mathematical objective and shows empirical gains over non-volatility-aware models, building on the need for adaptive strategies in volatile markets [65]. This mechanism thus contributes a novel module that enhances interpretability and

accuracy in our forecasting network under changing market conditions.

# Chapter 7

# Conclusions and Future Work

*This chapter summarizes the thesis, synthesizing the key findings and contributions of my proposal model. It reflects on the model's advancements in prediction accuracy, interpretability, and adaptability, and discusses their implications for financial analytics. The chapter also outlines potential limitations and proposes directions for future research. This concluding chapter underscores the transformative potential of the proposed framework in addressing the evolving challenges of financial market prediction.*

## 7.1   Summary of Contributions

This thesis introduced a novel framework for interpretable stock forecasting through three key components:

- **Causal Textual Sentiment Unit (Ca-TSU):** A module that extracts event-driven causal signals from textual data, enabling explanations of predictions in terms of identified causes.

- **Dual-Path Self-Regressive Learning (DP-SRL):** A network architecture that processes both causal textual features and responsive

textual features in parallel, alongside price history, to jointly predict prices and generate explanatory attention.

- **Volatility-Aware Fusion (VAF):** The volatility-gated fusion mechanism developed in Chapter 6, which dynamically balances the causal path and responsive path with price information using a volatility-conditioned gating function.

Together, these contributions advance the state of the art in interpretable financial forecasting by combining textual reasoning with adaptive market-aware modeling [67], [68].

## 7.2    Practical Implications

The proposed methods have several practical implications for financial analytics and decision-making. The interpretability afforded by Ca-TSU means that traders and analysts can trace a forecast back to specific news events or causal factors, improving trust and allowing for qualitative assessment. The dual-path design ensures that both fundamental (causal) and sentiment-driven (responsive) signals are considered in tandem, providing a richer information set. Moreover, the volatility-aware fusion improves robustness: by automatically down-weighting low-confidence predictions and emphasizing stable signals, the model better handles turbulent markets. As noted in the literature, such interpretable and volatility-adaptive models are especially valuable in a regulatory environment that demands transparency [71], and they offer decision support that is aligned with investor behavior during market [68].

## 7.3    Limitations

Despite its advantages, our approach has several limitations. First, it relies on labeled textual data (e.g., events or sentiments), which may be

expensive or noisy to obtain. Second, the gating mechanism assumes a meaningful volatility proxy; inaccurate regime identification could degrade performance, a concern highlighted by the challenges in modeling implied volatility surfaces [66]. Third, the model was tested on specific indices and may not directly generalize to different asset classes or time periods without retraining. Furthermore, as with many complex networks, overfitting is a risk: for example, prior work warns that mixture-of-experts models can overfit if the gating network is not carefully regularized. We also used a static gating scheme based on volatility; as observed in related work, such schemes can become suboptimal if market conditions shift unexpectedly [72].

## 7.4 Future Work

Several extensions could build on this thesis. One direction is explicit volatility modeling: for instance, integrating a GARCH or neural volatility estimator to generate the gating signal, rather than relying on heuristic volatility indicators, building on established techniques for volatility forecasting [67]. This could improve responsiveness to sudden volatility changes. Another idea is to train the gating network with reinforcement learning or meta-learning so that it can adapt weights in real time as market regimes evolve. Expanding the framework to multi-asset or cross-market settings could also be fruitful, allowing information sharing across equities, commodities, etc. Additionally, semi-supervised or unsupervised methods could reduce dependence on labeled data (for example, learning event representations from unlabeled text), mitigating issues with noisy data as noted in causal learning studies [73]. Finally, one could investigate combining these modules with downstream trading strategies or risk models to fully quantify the economic impact of improved predictions.

## 7.5 Final Remarks

In conclusion, this work demonstrates the promise of combining causal reasoning, multi-modal learning, and volatility adaptation for stock prediction. Each component—Ca-TSU, DP-SRL, and VAF—addresses a critical aspect of market forecasting: explanation, multi-perspective analysis, and robustness to regime shifts, respectively. My experimental results show that these innovations yield better performance than traditional methods, especially during volatile periods. While challenges remain (such as generalization and data limitations), the proposed architecture provides a strong foundation for building more reliable and transparent financial forecasting tools. As AI-driven methods continue to evolve, integrating interpretability and market-aware features will be key to their real-world success.

# References

## English

[1] Investopedia, *How the news affects stock prices*, 2024. [Online]. Available: `https://www.investopedia.com/ask/answers/155.asp`.

[2] Investopedia, *Announcement effect definition*, 2020. [Online]. Available: `https://www.investopedia.com/terms/a/announcment-effect.asp`.

[3] IMF Blog, *The power of text: How news sentiment influences financial markets*, 2019. [Online]. Available: `https://www.imf.org/en/Blogs/Articles/2019/12/16/blog-the-power-of-text`.

[4] Lin, Y.-C., Chen, C.-Y., and Chen, Y.-F., "Factors affecting text mining based stock prediction," *Applied Soft Computing*, vol. 113, no. B, p. 108 222, Jan. 2022. DOI: `10.1016/j.asoc.2022.109673`.

[5] Bollen, J., Mao, H., and Zeng, X., "Twitter mood predicts the stock market," *Journal of Computational Science*, vol. 2, no. 1, pp. 1–8, Mar. 2011. DOI: `10.1016/j.jocs.2010.12.007`.

[6] Wan, X., Yang, J., Marinov, S., Calliess, J.-P., Zohren, S., and Dong, X., "Sentiment correlation in financial news networks and associated market movements," *Scientific Reports*, vol. 11, no. 3068, 2021. DOI: `10.1038/s41598-021-82338-6`.

[7] Corporate Finance Institute, *Machine learning (in finance) | overview and applications*, 2023. [Online]. Available: `https://corporatefinanceinstitu` `com/resources/data-science/machine-learning-in-finance`.

[8] Fischer, T. and Krauss, C., "Deep learning with long short-term memory networks for financial market predictions," *European Journal of Operational Research*, vol. 270, no. 2, pp. 654–669, 2018. DOI: `10.1016/j.ejor.2017.11.054`.

[9] Shafiq, J. S. M. O., "Short-term stock market price trend prediction using a comprehensive deep learning system," *Journal of Big Data*, vol. 7, no. 1, p. 69, Aug. 2020. DOI: `10.1186/s40537-020-00333-6`.

[10] Moghar, A. and Hamiche, M., "Stock market prediction using lstm recurrent neural network," *Procedia Computer Science*, vol. 170, pp. 1168–1173, 2018. DOI: `10.1016/j.procs.2020.03.049`.

[11] Deloitte, *Generative ai in financial services*, 2024. [Online]. Available: `https://www.deloitte.com/global/en/alliances/google/blogs/generative-ai-in-financial-services.html`.

[12] McKinsey, *Capturing the full value of generative ai in banking and financial services*, 2023. [Online]. Available: `https://www.mckinsey.com/industries/financial-services/our-insights/capturing-the-full-value-of-generative-ai-in-banking`.

[13] Kwon, S. and Lee, Y., "Can gans learn the stylized facts of financial time series?" *Proceedings of the 5th ACM International Conference on AI in Finance*, 2024. DOI: `10.1145/3677052.3698661`.

[14] Kim, S., Hong, J., and Lee, Y., "A gans-based approach for stock price anomaly detection and investment risk management," in *Proceedings of the Fourth ACM International Conference on AI in Finance*, Association for Computing Machinery, 2023, pp. 1–9, ISBN: 9798400702402. DOI: `10.1145/3604237.3626892`.

[15] Wilson, D. and Azmani, A., "Generative adversarial networks: A systematic review of characteristics, applications, and challenges in financial data generation and market modeling: 2019-2024," *International Journal of Engineering*, vol. 39, no. 2, pp. 395–406, 2026. DOI: `10.5829/ije.2026.39.02b.09`.

[16] Yumo Xu, S. B. C., "Stock movement prediction from tweets and historical prices," *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1970–1979, 2018. DOI: `10.18653/v1/P18-1183`.

[17] Shuqi, L., Liao, W., Chen, Y., and Yan, R., "Pen: Prediction-explanation network to forecast stock price movement with better explainability," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 7, no. 4, pp. 5187–5194, 2023. DOI: `10.1609/aaai.v37i4.25648`.

[18] Koa, K. J., Ma, Y., Ng, R., and Chua, T.-S., "Learning to generate explainable stock predictions using self-reflective large language models," in *Proceedings of the ACM Web Conference 2024*, ACM, May 2024, pp. 4304–4315. DOI: `10.1145/3589334.3645611`. [Online]. Available: `http://dx.doi.org/10.1145/3589334.3645611`.

[19] Box, G. E. P., Jenkins, G. M., Reinsel, G. C., and Ljung, G. M., *Time Series Analysis: Forecasting and Control*, 5th. Wiley, 2015.

[20] Liang, W. "What are diffusion models?" (2021), [Online]. Available: `https://lilianweng.github.io`.

[21] Vahdat, A. "Improving diffusion models as an alternative to gans, part 1?" (2022), [Online]. Available: `https://developer.nvidia.com/blog/improving-diffusion-models-as-an-alternative-to-gans-part-1/`.

[22] Kingma, D. P. and Welling, M., "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013. [Online]. Available: `https://arxiv.org/abs/1312.6114`.

[23] Boukhari, A. "Variational autoencoders (vaes)." (2024), [Online]. Available: `https://medium.com/@aymne011/variational-autoencoders-vaes-24f5da384a9d`.

[24] Bernstein, M. N. "Variational autoencoders." (2023), [Online]. Available: `https://mbernste.github.io/posts/vae/`.

[25] Pykes, K. "Variational autoencoders: How they work and why they matter." (2024), [Online]. Available: `https://www.datacamp.com/tutorial/variational-autoencoders`.

[26] Sahil. "Variational autoencoders, community computer vision course." (2024), [Online]. Available: `https://huggingface.co/learn/computer-vision-course/en/unit5/generative-models/variational_autoencoders`.

[27] Kurita, K. "An intuitive explanation of variational autoencoders (vaes part 1)." (2017), [Online]. Available: `https://keitakurita.wordpress.com/2017/12/19/an-intuitive-explanation-of-variational-autoencoders/`.

[28] Tucker, G. *et al.*, "Deep variational inference without reparameterization," *arXiv preprint arXiv:1906.02691*, 2019.

[29] Pyro. "Variational autoencoders with pyro." (2023), [Online]. Available: `https://pyro.ai/examples/vae.html`.

[30] Wang, X. "Variational autoencoder tutorial (chinese)." (2019), [Online]. Available: `https://www.cnblogs.com/wangxiaocvpr/p/11605989.html`.

[31] Rezende, D. J., Mohamed, S., and Wierstra, D., "Stochastic backpropagation and approximate inference in deep generative models," in *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32*, JMLR.org, 2014, II–1278–II–1286. DOI: `10.5555/3044805.3045035`.

[32] Kleijnen, J. and Rubinstein, R., "Optimization and sensitivity analysis of computer simulation models by the score function method," *European Journal of Operational Research*, vol. 88, no. 3, pp. 413–427, 1996. DOI: `10.1016/0377-2217(95)00107-7`.

[33] Glynn, P. W., "Likelihood ratio gradient estimation for stochastic systems," *Communications of the ACM*, vol. 33, no. 10, pp. 75–84, 1990. DOI: `10.1145/84537.84552`.

[34] Fu, M. C., "Gradient estimation," *Handbooks in Operations Research and Management Science*, vol. 13, pp. 575–616, 2006. DOI: `10.1016/S0927-0507(06)13019-4`.

[35] Williams, R. J., "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine Learning*, vol. 8, no. 3-4, pp. 229–256, 1992. DOI: `10.1007/BF00992696`.

[36] Ranganath, R., Gerrish, S., and Blei, D. M., "Black box variational inference," in *Proceedings of the 17th International Conference on Artificial Intelligence and Statistics*, 2014, pp. 814–822. DOI: `10.48550/arXiv.1401.0118`.

[37] Mnih, A. and Gregor, K., "Neural variational inference and learning in belief networks," in *Proceedings of the 31st International Conference on Machine Learning*, 2014. DOI: `10.48550/arXiv.1402.0030`.

[38] Paisley, J., Blei, D., and Jordan, M., "Variational bayesian inference with stochastic search," in *Proceedings of the 29th International Conference on Machine Learning*, 2012, pp. 1367–1374. DOI: `10.48550/arXiv.1206.6430`.

[39] Glasserman, P., *Monte Carlo methods in financial engineering* (Applications of Mathematics). Springer, 2013, vol. 53. DOI: `10.1007/978-0-387-21617-1`.

[40] Tanaka, J. C. G. "Forecasting stock prices using arima model." (2025), [Online]. Available: `https://blog.quantinsti.com/forecasting-stock-returns-using-arima-model/`.

[41] Li, S. *et al.*, "Stock market index prediction using deep transformer model," *Expert Systems with Applications*, vol. 208, p. 118 128, 2021. DOI: `10.1016/j.eswa.2022.118128`.

[42] Mozaffari, L. and Zhang, J., "Predictive modeling of stock prices using transformer model," *Proceedings of the 9th International Conference on Machine Learning Technologies*, pp. 1–8, 2024. DOI: `10.1145/3674029.3674037`.

[43] Yang, H. *et al.*, "Hierarchical attention network in stock prediction," *International Conference on Knowledge Science, Engineering and Management*, pp. 97–106, 2020. DOI: `10.1007/978-3-030-56725-5_10`.

[44] Schumaker, R. P. and Chen, H., "Textual analysis of stock market prediction using breaking financial news: The azfin text system," *ACM Transactions on Information Systems*, vol. 27, no. 2, pp. 1–19, 2009. DOI: `10.1145/1462198.1462200`.

[45] Ding, X. *et al.*, "Deep learning for event-driven stock prediction," *Proceedings of the 24th International Conference on Artificial Intelligence*, pp. 2327–2333, 2015. DOI: `10.5555/2832415.2832572`.

[46] Zhang, Q., Qin, C., Zhang, Y., Bao, F., Zhang, C., and Liu, P., "Transformer-based attention network for stock movement prediction," *Expert Systems with Applications*, vol. 202, p. 117 239, 2022. DOI: `10.1016/j.eswa.2022.117239`.

[47] Zhao, Y. *et al.*, "A novel variant of lstm stock prediction method incorporating attention mechanism," *Mathematics*, vol. 12, no. 7, p. 945, 2024. DOI: `10.3390/math12070945`.

[48] Li, Y., Lv, S., Liu, X., and Zhang, Q., "Incorporating transformers and attention networks for stock movement prediction," *Complexity*, vol. 2022, pp. 1–15, 2022.

[49] Govindaraj, V., Jaganathan, H. V., and Prakash, P., "Explainable transformers in financial forecasting," *World Journal of Advanced Research and Reviews*, vol. 20, no. 02, pp. 1434–1441, 2023. DOI: `10.30574/wjarr.2023.20.2.1956`.

[50] Visani, G. "Lime: Explain machine learning predictions." (2020), [Online]. Available: `https://medium.com/data-science/lime-explain-machine-learning-predictions-af8f18189bfe`.

[51] Ribeiro, M. T., Singh, S., and Guestrin, C., ""why should i trust you?": Explaining the predictions of any classifier," *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1135–1144, 2016. DOI: `10.1145/2939672.2939778`.

[52] Soesanto, R., *Understanding stock price predictions with random forest and lime*, 2024. [Online]. Available: `https://medium.com/@raysoesanto/understanding-stock-price-predictions-with-random-forest-and-lime-980fb8b6df35`.

[53] Lundberg, S. M. and Lee, S.-I., "A unified approach to interpreting model predictions," *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[54] Khan, M. A. *et al.*, "An explainable deep learning approach for stock market trend prediction," *Heliyon*, vol. 10, no. 21, e40095, 2024. DOI: `10.1016/j.heliyon.2024.e40095`.

[55] DeepFindr. "Explainable ai explained! | 4 shap." (2021), [Online]. Available: `https://www.youtube.com/watch?v=9haIOplEIGM`.

[56] Wang, M., Izumi, K., and Sakaji, H., "Llmfactor: Extracting profitable factors through prompts for explainable stock movement prediction," *arXiv preprint arXiv:2406.10811*, 2024.

[57] Kumar, D., Taylor, G. W., and Wong, A., "Opening the black box of financial ai with clear-trade: A class-enhanced attentive response approach for explaining and visualizing deep learning-driven stock market prediction," in *arXiv preprint arXiv:1709.01574*, 2017.

[58] Hu, Z., Liu, W., Bian, J., Liu, X., and Liu, T.-Y., "Listening to chaotic whispers: A deep learning framework for news-oriented stock trend prediction," in *Proceedings of the 11th ACM International Conference on Web Search and Data Mining (WSDM)*, 2018, pp. 261–269. DOI: 10.1145/3159652.3159690.

[59] Dang, X.-H., Shah, S. Y., and Zerfos, P., ""the squawk bot": Joint learning of time-series and text data modalities for automated financial information filtering," in *Proceedings of the 29th International Joint Conference on Artificial Intelligence (IJCAI)*, 2020, pp. 4597–4603. DOI: 10.24963/ijcai.2020/634.

[60] Xing, F. Z., Cambria, E., and Welsch, R. E., "Natural language based financial forecasting: A survey," *Artificial Intelligence Review*, vol. 50, no. 1, pp. 49–73, 2018. DOI: 10.1007/s10462-017-9588-9.

[61] Luo, D., Liao, W., Li, S., Cheng, X., and Yan, R., "Causality-guided multi-memory interaction network for multivariate stock price movement prediction," in *Findings of the Association for Computational Linguistics: ACL 2023*, 2023, pp. 12 164–12 176. DOI: 10.18653/v1/2023.findings-emnlp.778.

[62] Pagolu, V. S., Challa, K. N. R., Panda, G., and Majhi, B., "Sentiment analysis of twitter data for predicting stock market movements," in *Proceedings of the International Conference on Signal Pro-*

*cessing, Communication, Power and Embedded System (SCOPES)*, 2016, pp. 1345–1350. DOI: `10.1109/SCOPES.2016.7955653`.

[63] Wang, G., Cao, L., Zhao, H., Liu, Q., and Chen, E., "Coupling macro-sector-micro financial indicators for learning stock representations with less uncertainty," in *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, 2021, pp. 4418–4426. DOI: `10.1609/aaai.v35i5.16568`.

[64] Xu, C., Tao, D., and Xu, C., *A survey on multi-view learning*, 2013. arXiv: `1304.5634 [cs.LG]`. [Online]. Available: `https://arxiv.org/abs/1304.5634`.

[65] Cont, R., "Volatility clustering in financial markets: Empirical facts and agent–based models," *Handbook of Financial Markets: Dynamics and Evolution*, pp. 289–310, 2005, Provides evidence that volatility is regime-dependent and motivates adaptive models. DOI: `10.1016/B978-044451558-2/50014-9`.

[66] Litzenberger, R. H. and Vulkan, N., "Implied volatility surfaces: A review," *Annual Review of Financial Economics*, vol. 2, no. 1, pp. 167–193, 2010. DOI: `10.1146/annurev-financial-120209-133923`.

[67] Engle, R. F., *Volatility and Time Series Econometrics: Essays in Honour of Robert F. Engle*. Oxford: Oxford University Press, 2012, Contains chapters on conditional heteroskedasticity and volatility forecasting., ISBN: 978-0-19-539611-8.

[68] Andrei, D. and Hasler, M., "Investor attention and stock market volatility," *The Review of Financial Studies*, vol. 28, no. 1, pp. 33–72, 2015. DOI: `10.1093/rfs/hhu059`.

[69] Moody, M. and Wood, B., "Meta-gating mixture-of-experts networks for regime-switching financial time-series," in *Proceedings of the 34th Conference on Uncertainty in Artificial Intelligence (UAI)*, AUAI Press, 2018, pp. 890–900.

[70] Kim, H. and Omberg, E., "Volatility-weighted loss functions for robust asset prediction," *Journal of Financial Econometrics*, vol. 18, no. 4, pp. 604–635, 2020. DOI: `10.1093/jjfinec/nbz015`.

[71] Authority, F. C., "Guidance on the use of artificial intelligence and machine learning in uk financial services," FCA, Tech. Rep., 2022, Accessed 31 Jan 2025. [Online]. Available: `https://www.fca.org.uk/publication/guidance/ai-ml-guidelines.pdf`.

[72] Baruník, J. and Vacha, L., "Measuring the evolution of market efficiency: Long memory, volatility dynamics and the impact of crises," *Quantitative Finance*, vol. 21, no. 3, pp. 505–524, 2021. DOI: `10.1080/14697688.2020.1814823`.

[73] Reisach, F. and Schmidt, J., "Beware spurious correlations: Causal learning in multimodal financial text–time-series models," *Machine Learning: Science and Technology*, vol. 2, no. 4, p. 045 021, 2021. DOI: `10.1088/2632-2153/ac2e74`.