

**ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH**  
**TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN**



**UIT**  
**TRƯỜNG ĐẠI HỌC**  
**CÔNG NGHỆ THÔNG TIN**

**BÁO CÁO ĐỒ ÁN MÔN HỌC**  
**CÔNG NGHỆ DỮ LIỆU LỚN**

**ĐỀ TÀI:**

**Hệ thống truy xuất và dự đoán dữ liệu tài chính theo thời gian  
thực từ sàn giao dịch điện tử Coinbase**

***Lớp:*** IE212.P22

***GVHD:*** TS. Hà Minh Tân

***Nhóm sinh viên thực hiện:***

Nguyễn Đặng Minh Quan 22521183

Trần Mạnh Phúc 22521142

An Nhất Lâm 22520732

*Thành phố Hồ Chí Minh, tháng 5 năm 2025.*

## NHẬN XÉT CỦA GIÁO VIÊN HƯỚNG DẪN

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

....., ngày... tháng.....năm 2025

Người nhận xét  
(Ký tên và ghi rõ họ tên)

## LỜI CẢM ƠN

Trước hết, nhóm chúng em xin được gửi lời cảm ơn chân thành và sâu sắc nhất tới Tiến sĩ Hà Minh Tân – giảng viên môn Công nghệ Dữ liệu lớn. Chúng em vô cùng trân trọng và biết ơn Thầy vì đã luôn tận tâm truyền đạt những kiến thức quý báu, giúp chúng em có được nền tảng vững chắc để thực hiện đồ án này.

Trong suốt quá trình học tập và làm việc với đề tài, Thầy không chỉ cung cấp những bài giảng súc tích, dễ hiểu mà còn luôn sẵn sàng lắng nghe, hướng dẫn và góp ý tận tình. Chính sự hỗ trợ và định hướng của Thầy đã giúp chúng em vượt qua nhiều khó khăn, từng bước hoàn thiện nội dung nghiên cứu một cách hiệu quả nhất.

Mặc dù nhóm đã cố gắng hết sức để áp dụng kiến thức đã học và tìm hiểu thêm từ các nguồn tài liệu khác, nhưng với giới hạn về kinh nghiệm và thời gian, chắc chắn đồ án vẫn còn những điểm chưa hoàn hảo. Chúng em rất mong nhận được những nhận xét và góp ý từ Thầy để có thể học hỏi thêm, hoàn thiện kỹ năng cho những nghiên cứu tiếp theo.

Nhóm xin chân thành cảm ơn Thầy TS. Hà Minh Tân đã tận tình hướng dẫn và hỗ trợ em trong suốt quá trình thực hiện đồ án.

# MỤC LỤC

<b>Chương 1. TỔNG QUAN ĐỀ TÀI .....</b>	<b>7</b>
1.1. Tính cấp thiết của đề tài .....	7
1.2. Mục tiêu nghiên cứu.....	8
1.3. Đối tượng và phạm vi nghiên cứu.....	9
1.4. Nội dung đề tài.....	10
<b>Chương 2. CƠ SỞ KIẾN TRÚC XỬ LÝ THỜI GIAN THỰC .....</b>	<b>11</b>
2.1. Bài toán xử lý thời gian thực .....	11
2.1.1. Định nghĩa, phân loại, yêu cầu .....	11
2.1.2. Cấu trúc kappa trong xây dựng hệ thống xử lý dữ liệu thời gian thực .....	12
2.2. Apache Kafka - Xương sống của streaming platform.....	14
2.3. Apache Spark.....	15
2.4. Apache Cassandra .....	16
2.5. Thiết kế pipeline dữ liệu.....	17
2.5.1. Thu thập dữ liệu - Coinbase WebSocket API .....	18
2.5.2. Kafka Producer Service – Hàm truyền dữ liệu tới Kafka .....	18
2.5.3. Apache Kafka - xương sống của streaming.....	19
2.5.4. Lưu trữ dữ liệu thô .....	20
2.5.5. Xử lý dữ liệu thời gian thực .....	21
2.5.6. Apache Cassandra - Lưu trữ và truy xuất dữ liệu .....	21
2.5.7. Dự đoán và trực quan hóa .....	22
2.5.8. Triển khai hệ thống: .....	24
<b>Chương 3. CÁC CÔNG VIỆC LIÊN QUAN .....</b>	<b>25</b>
3.1. Các phương pháp dự đoán giá tiền điện tử.....	25
3.2. Pipeline dữ liệu trong dự đoán tài chính.....	27
3.3. Khoảng trống nghiên cứu.....	28
<b>Chương 4. PHƯƠNG PHÁP ÁP DỤNG CHO BÀI TOÁN.....</b>	<b>29</b>
4.1. Mô tả dữ liệu.....	29
4.2. Chuẩn hóa dữ liệu .....	32
4.2.1. Min-Max Scaling.....	33
4.2.2. Robust Scaling.....	33
4.2.3. Phân nhóm đặc trưng theo phương pháp chuẩn hóa .....	35
4.3. Tăng cường dữ liệu (Data Augmentation) .....	35
4.3.1. Che cục bộ (Local Mean Masking).....	36
4.3.2. Thêm nhiễu Gaussian (Adaptive Gaussian Noise) .....	36
4.3.3. Tỷ lệ ngẫu nhiên (Smart Scaling).....	36

4.3.4. Giãn thời gian (Time Warping).....	37
4.3.5. Bỏ đặc trưng ngẫu nhiên (Feature Dropout).....	37
4.4. Các Mô hình Dự đoán: .....	38
4.4.1. Mạng LSTM thuần túy.....	38
4.4.2. Mô hình CNN-LSTM.....	38
4.4.3. Mô hình CNN-LSTM Tích hợp Attention .....	39
<b>Chương 5. KẾT QUẢ THỰC NGHIỆM VÀ THẢO LUẬN .....</b>	<b>40</b>
5.1. Các Mô hình Dự đoán: .....	40
5.1.1. Dữ liệu và cấu trúc chuỗi đầu vào .....	40
5.1.2. Mô hình huấn luyện.....	41
a) Mô hình LSTM thuần túy .....	41
b) Mô hình CNN-LSTM.....	41
c) Mô hình CNN-LSTM Attention .....	42
5.1.3. Quy trình huấn luyện .....	43
5.1.4. Phân chia dữ liệu .....	44
5.1.5. Môi trường thực nghiệm .....	44
5.2. Tiêu chí đánh giá .....	45
5.3. Phân tích kết quả thực nghiệm.....	46
5.3.1 <i>Hiệu suất mô hình LSTM</i> .....	47
5.3.2. <i>Hiệu suất mô hình CNN-LSTM</i> .....	51
5.3.3. <i>Hiệu suất mô hình CNN-LSTM kết hợp Attention</i> .....	55
<b>Chương 6. KẾT LUẬN – ĐỀ XUẤT .....</b>	<b>59</b>
6.1. Kết luận .....	59
6.2. Đề xuất.....	60
<b>TÀI LIỆU THAM KHẢO.....</b>	<b>63</b>

## DANH MỤC HÌNH ẢNH

Hình 2.1. Phác thảo cơ bản của kiến trúc Kappa .....	13
Hình 2.2. Kiến trúc cơ bản của Apache Kafka .....	14
Hình 2.3. Apache Spark .....	15
Hình 2.4. Apache Cassandra .....	16
Hình 2.5. Pipeline dữ liệu thời gian thực .....	17
Hình 2.6. Visualize giá tiền điện tử theo thời gian thực .....	22
Hình 2.7. Visualize giá tiền điện tử dự đoán theo thời gian thực .....	23
Hình 3.1. Đánh giá trianing test .....	25
Hình 3.2. Đánh giá training error .....	25
Hình 3.3. Cấu trúc cơ bản của bài báo Nazzer et al., 2017 .....	27
Hình 5.1. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình LSTM – Mẫu 1. ....	48
Hình 5.2. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình LSTM – Mẫu 2. ....	48
Hình 5.3. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình LSTM – Mẫu 3. ....	49
Hình 5.4. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình LSTM – Mẫu 4. ....	49
Hình 5.5. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình LSTM – Mẫu 5. ....	50
Hình 5.6. Phân phối lỗi của mô hình LSTM.....	50
Hình 5.7. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM – Mẫu 1. ....	52
Hình 5.8. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM – Mẫu 2. ....	52
Hình 5.9. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM – Mẫu 3. ....	53
Hình 5.10. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM – Mẫu 4. ....	53
Hình 5.11. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM – Mẫu 5. ....	54
Hình 5.12. Phân phối lỗi của mô hình CNN-LSTM.....	54
Hình 5.13. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM -Attention– Mẫu 1.....	56
Hình 5.14. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM -Attention– Mẫu 2.....	56
Hình 5.15. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM -Attention– Mẫu 3.....	57
Hình 5.16. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM -Attention– Mẫu 4.....	57
Hình 5.17. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM -Attention– Mẫu 5.....	58
Hình 5.18. Phân phối lỗi của mô hình CNN-LSTM-Attention. ....	58

## **Chương 1. TỔNG QUAN ĐỀ TÀI**

### **1.1. Tính cấp thiết của đề tài**

Cùng với sự phát triển mạnh mẽ của công nghệ số và thị trường tài chính điện tử, việc truy xuất và dự đoán dữ liệu tài chính theo thời gian thực đã trở nên cấp thiết để hỗ trợ các nhà đầu tư, nhà giao dịch và các tổ chức tài chính đưa ra quyết định chính xác và kịp thời. Các sàn giao dịch điện tử như Coinbase cung cấp khối lượng dữ liệu khổng lồ về giá cả, khối lượng giao dịch và xu hướng thị trường, tạo cơ hội cho các cá nhân và tổ chức tiếp cận thông tin tài chính một cách nhanh chóng và hiệu quả. Tuy nhiên, để khai thác tối đa tiềm năng của dữ liệu này, cần có các hệ thống thông minh có khả năng xử lý và phân tích dữ liệu theo thời gian thực.

Việc truy xuất và dự đoán dữ liệu tài chính theo thời gian thực không chỉ đòi hỏi công nghệ tiên tiến mà còn yêu cầu người dùng có khả năng phân tích và đưa ra quyết định dựa trên thông tin được cung cấp. Với khối lượng dữ liệu lớn và sự biến động không ngừng của thị trường, việc xử lý thủ công trở nên bất khả thi, dẫn đến nhu cầu về các hệ thống tự động hóa có khả năng thu thập, phân tích và đưa ra dự đoán chính xác. Tuy nhiên, do tính chất phức tạp và tốc độ thay đổi nhanh của dữ liệu tài chính, nhiều hệ thống hiện tại gặp khó khăn trong việc cung cấp thông tin đáng tin cậy hoặc cập nhật kịp thời, khiến người dùng bỏ lỡ cơ hội hoặc đối mặt với rủi ro tài chính. Nhờ có hệ thống này, người dùng không chỉ cải thiện hiệu suất giao dịch, đưa ra quyết định nhanh chóng và chính xác hơn mà còn giảm thiểu rủi ro trong môi trường thị trường biến động.

Chính vì những lý do trên, nhóm chọn đề tài: “Xây dựng hệ thống truy xuất và dự đoán dữ liệu tài chính theo thời gian thực từ sàn giao dịch điện tử Coinbase”. Nhóm hy vọng rằng đồ án sẽ giải quyết các vấn đề cấp bách và quan trọng liên quan đến việc xử lý và khai thác dữ liệu tài chính, từ đó mang lại giá trị thiết thực cho người dùng trong bối cảnh thị trường tài chính số ngày càng phát triển.

## 1.2. Mục tiêu nghiên cứu

**Mục tiêu tổng quát:** Xây dựng một hệ thống pipeline dữ liệu thời gian thực từ sàn giao dịch điện tử Coinbase để truy xuất, xử lý và dự đoán dữ liệu tài chính, từ đó:

- Cung cấp các dự đoán giá tiền điện tử chính xác và cá nhân hóa cho người dùng (nhà đầu tư, nhà giao dịch).
- Tăng cường hiệu quả ra quyết định đầu tư thông qua các dự đoán giá dựa trên dữ liệu thời gian thực.
- Giảm thiểu rủi ro tài chính bằng cách cung cấp thông tin thị trường kịp thời và đáng tin cậy.
- Hỗ trợ quản lý chiến lược giao dịch và tối ưu hóa danh mục đầu tư thông qua các công cụ trực quan hóa và phân tích dữ liệu.

**Mục tiêu cụ thể của đề tài bao gồm:**

- **Thu thập và xử lý dữ liệu:** Thiết lập pipeline dữ liệu thời gian thực sử dụng Coinbase WebSocket API để thu thập dữ liệu thị trường tiền điện tử, kết hợp Kafka Producer và Kafka Broker để truyền dữ liệu, và PySpark Structured Streaming để xử lý dữ liệu theo thời gian thực.
- **Lưu trữ dữ liệu:** Chuẩn hóa và lưu trữ dữ liệu thô vào AWS S3 (MinIO), đồng thời sử dụng Apache Cassandra để lưu trữ dữ liệu lịch sử dài hạn, đảm bảo hiệu suất đọc/ghi cao và khả năng mở rộng.
- **Huấn luyện và triển khai mô hình dự đoán:** Áp dụng các mô hình học sâu như LSTM, CNN-LSTM và CNN-LSTM Attention để huấn luyện và dự đoán giá tiền điện tử, sử dụng dữ liệu từ Cassandra để cung cấp dự đoán thời gian thực.
- **Trực quan hóa dữ liệu:** Xây dựng giao diện trực quan hóa bằng Grafana để hiển thị dữ liệu thị trường và kết quả dự đoán theo thời gian thực, hỗ trợ người dùng theo dõi xu hướng và đưa ra quyết định.
- **Tích hợp pipeline:** Đảm bảo luồng dữ liệu từ Coinbase qua các giai đoạn xử lý, lưu trữ, dự đoán và trực quan hóa hoạt động liền mạch, tối ưu hóa độ trễ và độ chính xác của dự đoán.



### 1.3. Đối tượng và phạm vi nghiên cứu

Đối tượng nghiên cứu của đề tài là bài toán truy xuất và dự đoán giá tiền điện tử theo thời gian thực từ sàn giao dịch Coinbase, sử dụng dữ liệu thị trường bao gồm thông tin giao dịch, giá cả và xu hướng thị trường. Bài toán được mô tả cụ thể với đầu vào và đầu ra như sau:

- **Input:** Dữ liệu lớn từ sàn giao dịch Coinbase, bao gồm:
  - Thông tin về các giao dịch tiền điện tử (giá, khối lượng).
  - Dữ liệu giá lịch sử và thời gian thực.
  - Chuỗi dữ liệu thị trường theo thời gian (biến động giá, xu hướng).
- **Output:** Dự đoán top  $k$  ( $k \in \mathbb{N}^*$ ) giá tiền điện tử trong tương lai gần hoặc các khuyến nghị giao dịch phù hợp nhất cho người dùng.

Phạm vi nghiên cứu: Dữ liệu đầu vào của bài toán sẽ được thu thập từ Coinbase WebSocket API, bao gồm thông tin giao dịch và giá cả tiền điện tử theo thời gian thực, kết hợp với dữ liệu lịch sử từ các nguồn bổ sung nếu cần. Sử dụng các công nghệ lưu trữ và xử lý dữ liệu lớn như:

- **Apache Kafka** để truyền dữ liệu thời gian thực.
- **PySpark Structured Streaming** để xử lý dữ liệu lớn theo thời gian thực.
- **AWS S3 (MinIO)** để lưu trữ dữ liệu thô.
- **Apache Cassandra** để lưu trữ dữ liệu lịch sử dài hạn với hiệu suất cao.
- Xây dựng các mô hình dự đoán giá tiền điện tử sử dụng PySpark tích hợp với các mô hình học sâu (LSTM, CNN-LSTM, CNN-LSTM Attention) được triển khai bằng PyTorch.

Cuối cùng, triển khai một giao diện trực quan hóa bằng Grafana để hiển thị dữ liệu thị trường và kết quả dự đoán theo thời gian thực, hỗ trợ người dùng theo dõi và ra quyết định.

Nghiên cứu này tập trung vào phát triển các giải pháp kỹ thuật để dự đoán xu hướng giá tài sản số từ dữ liệu thời gian thực, mà không triển khai giao dịch thực tế hoặc cung cấp tư vấn đầu tư. Phạm vi dữ liệu dự báo được giới hạn ở một số cặp giao dịch cụ thể, với ưu tiên cho cặp BTC-USD, nhằm đảm bảo tính khả thi và trọng tâm của đề tài. Tổng thể, đề tài mang tính ứng dụng cao, là tiền đề quan trọng cho việc xây dựng các hệ thống giao dịch tài chính tự động trong tương lai hoặc cung cấp công cụ phân tích cho nhà đầu tư cá nhân.

## 1.4. Nội dung đề tài

**Chương 1: Tổng quan đề tài** – Chương này giới thiệu tính cấp thiết của việc dự đoán giá tiền điện tử từ Coinbase để hỗ trợ nhà đầu tư, đặt mục tiêu xây dựng hệ thống truy xuất và dự đoán thời gian thực, xác định đối tượng là dữ liệu giao dịch và phạm vi sử dụng các công nghệ như Kafka, PySpark, Cassandra, cùng phương pháp tiếp cận từ thu thập đến trực quan hóa dữ liệu.

**Chương 2: Cơ sở lý thuyết** – Chương trình bày nền tảng lý thuyết về pipeline dữ liệu thời gian thực, xử lý dữ liệu lớn với Kafka, PySpark, Cassandra, và lưu trữ trên S3/MinIO, cùng các mô hình học sâu (LSTM, CNN-LSTM, CNN-LSTM Attention) để dự đoán chuỗi thời gian, kết hợp vai trò của Grafana trong trực quan hóa.

**Chương 3: Phương pháp thực hiện** – Chương mô tả phương pháp xây dựng hệ thống, từ thu thập dữ liệu qua Coinbase WebSocket API, truyền dữ liệu bằng Kafka, xử lý với PySpark, lưu trữ trên S3/MinIO và Cassandra, huấn luyện mô hình LSTM, CNN-LSTM, CNN-LSTM Attention bằng PyTorch, đến hiển thị kết quả trên Grafana.

**Chương 4: Kết quả và thảo luận** – Chương thống kê kết quả thực nghiệm của các mô hình dự đoán giá tiền điện tử, so sánh hiệu suất trên dữ liệu Coinbase, thảo luận ưu nhược điểm và độ chính xác của từng mô hình, đánh giá hiệu quả của pipeline dữ liệu và giao diện Grafana.

**Chương 5: Kết luận – đề xuất** – Chương cuối cùng cũng là chương sẽ khái quát toàn bộ những phương pháp đề xuất, đề cập các kết quả, ưu điểm mà phương pháp chúng tôi đề xuất. Sau cùng chúng tôi muốn đưa ra các hướng nghiên cứu, đề xuất phát triển đề tài trong tương lai có liên quan đến đề tài hiện tại chúng tôi thực hiện.

## Chương 2. CƠ SỞ KIẾN TRÚC XỬ LÝ THỜI GIAN THỰC

### 2.1. Bài toán xử lý thời gian thực

#### 2.1.1. Định nghĩa, phân loại, yêu cầu

Xử lý dữ liệu thời gian thực (real-time data processing) là quá trình thu thập, xử lý và phân tích dữ liệu ngay khi chúng được tạo ra, với độ trễ tối thiểu để đảm bảo tính kịp thời của thông tin [6]. Theo định nghĩa của Stonebraker et al. (2005), một hệ thống xử lý thời gian thực phải đáp ứng ba tiêu chí cơ bản: (i) xử lý dữ liệu liên tục (continuous data streams), (ii) đưa ra kết quả với độ trễ thấp, và (iii) có khả năng xử lý khối lượng dữ liệu lớn với tốc độ cao [7].

Trong bối cảnh thị trường tiền điện tử, xử lý dữ liệu thời gian thực đóng vai trò then chốt do một số đặc điểm riêng biệt của thị trường này. Thứ nhất, thị trường tiền điện tử hoạt động 24/7 không có thời gian đóng cửa như thị trường chứng khoán truyền thống, tạo ra luồng dữ liệu liên tục với khối lượng khổng lồ [8]. Thứ hai, tính biến động cao (high volatility) của giá tiền điện tử đòi hỏi các hệ thống phải có khả năng nắm bắt và phản ứng với mọi biến động giá trong thời gian ngắn nhất [9].

Nghiên cứu của Gandai và Halaburda (2016) đã chỉ ra rằng độ trễ trong việc cập nhật giá có thể dẫn đến những quyết định giao dịch sai lầm và tổn thất đáng kể trong thị trường tiền điện tử [10]. Tương tự, Dyhrberg (2016) nhấn mạnh rằng khả năng xử lý dữ liệu real-time là yếu tố quyết định trong việc xây dựng các mô hình dự đoán giá hiệu quả cho Bitcoin và các cryptocurrency khác [11].

Xử lý dữ liệu thời gian thực có thể được phân loại dựa trên mức độ nghiêm ngặt của yêu cầu thời gian, đặc biệt trong các ứng dụng tài chính:

1. **Hard real-time:** Yêu cầu độ trễ cố định, thường dưới 1ms, áp dụng trong các hệ thống giao dịch tần suất cao (high-frequency trading - HFT) nơi mỗi microgiây đều quan trọng. Ví dụ, các thuật toán HFT sử dụng dữ liệu thời gian thực để đặt lệnh mua/bán trong tích tắc [9].

2. **Soft real-time:** Chấp nhận độ trễ biến động trong khoảng cho phép, thường từ 1ms đến 1s. Điều này phù hợp với các ứng dụng như giám sát giá tiền điện tử trên sàn Coinbase, nơi độ trễ nhỏ không ảnh hưởng nghiêm trọng đến kết quả.
3. **Near real-time:** Độ trễ dao động từ 1s đến vài phút, phù hợp cho các tác vụ phân tích hoặc dự báo dài hạn hơn, chẳng hạn như dự đoán xu hướng giá trong vài giờ tới.

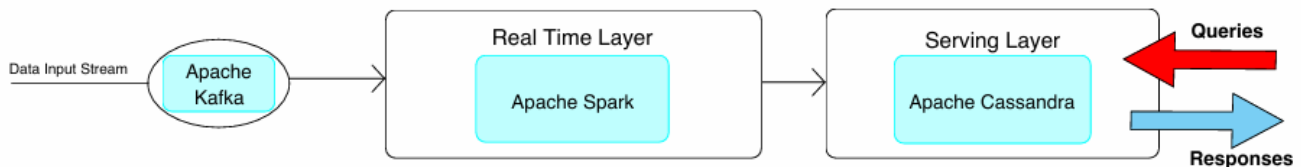
Theo đó, một hệ thống xử lý dữ liệu thời gian thực sẽ bao gồm những yêu cầu và những thách thức:

1. **Độ trễ:** Độ trễ trong hệ thống xử lý thời gian thực là khoảng thời gian từ khi một sự kiện giao dịch xảy ra đến khi kết quả xử lý (như dự đoán giá) được tạo ra [12]. Trong dự đoán giá tiền điện tử, độ trễ tổng thể bao gồm các giai đoạn thu thập dữ liệu (10-50ms), xử lý dữ liệu (50-200ms), thực hiện dự đoán bằng mô hình học sâu (100-500ms), và lưu trữ kết quả (10-30ms). Hirsche (2021) nhấn mạnh rằng độ trễ end-to-end dưới 1 giây là yêu cầu thiết yếu để đảm bảo tính cạnh tranh trong thị trường tiền điện tử [13].
2. **Thông lượng:** Thông lượng đo lường số lượng giao dịch hoặc sự kiện được xử lý mỗi giây. Trong thị trường tiền điện tử, thông lượng cần đạt 1,000-5,000 messages/giây trong điều kiện bình thường, tăng lên 10,000-50,000 trong giờ cao điểm, và có thể vượt 100,000 khi thị trường biến động mạnh.
3. **Tính nhất quán:** Tính nhất quán dữ liệu trong hệ thống phân tán là một thách thức lớn, như định lý CAP của Brewer (2000) nêu: không thể đồng thời đảm bảo nhất quán, khả dụng, và chịu lỗi phân vùng [15]. Để ưu tiên hiệu năng trong dự đoán giá tiền điện tử, hệ thống thường chấp nhận mô hình eventual consistency, sử dụng timestamp để sắp xếp dữ liệu và các kỹ thuật như window-based processing để xử lý dữ liệu đến muộn.
4. **Khả năng chịu lỗi là mở rộng:** Hệ thống cần mở rộng để xử lý khối lượng dữ liệu tăng đột biến. Horizontal scaling (thêm node, phân chia dữ liệu, cân bằng tải) và vertical scaling (tăng tài nguyên, tối ưu thuật toán, caching) là hai hướng chính. Thêm vào đó, nếu có lỗi xảy ra, hệ thống vẫn phải đảm bảo hoạt động thông qua các cơ chế như replication (nhân bản dữ liệu hoặc checkpointing (lưu trạng thái định kỳ)).

### 2.1.2. Cấu trúc kappa trong xây dựng hệ thống xử lý dữ liệu thời gian thực

Kiến trúc Lambda, được giới thiệu bởi Nathan Marz (2011), đã đang là giải pháp phổ biến cho việc xử lý dữ liệu lớn với cả batch và stream processing. Việc sử dụng kiến trúc Lambda cho phép các kỹ sư xây dựng luồng xử lý dữ liệu luồng để phục vụ nhu cầu truy vấn dữ liệu thời gian thực. Tuy nhiên, nó cũng yêu cầu duy trì hai cơ sở mã khác nhau, một cho batch và một cho thời gian thực. Điều đó cũng gây ra khó khăn khi phải xử lý logic nghiệp vụ giữa các mã nguồn xử lý real-time và batch sao cho phù hợp.

Để khắc phục những hạn chế này, người đồng sáng lập của Apache Kafka, Jay Kreps đã đề xuất sử dụng kiến trúc Kappa cho các hệ thống xử lý luồng. Ý tưởng chính của Kreps là lưu các dữ liệu từ nguồn dữ liệu có cấu trúc vào Kafka, chẳng hạn như bảng Hive của Apache. Sau đó, chỉ cần chạy lại các xử lý luồng trên các topic Kafka này, nhờ đó mã nguồn được thống nhất giữa cả luồng xử lý stream và batch.



Hình 2.1. Phác thảo cơ bản của kiến trúc Kappa

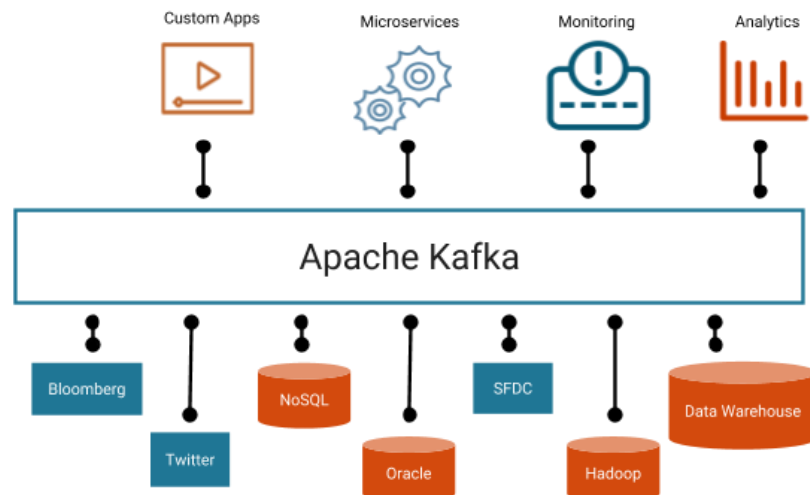
Hình 2.1 cung cấp phác thảo cơ bản của kiến trúc Kappa. Một lần nữa, kiến trúc này yêu cầu một luồng dữ liệu. Đầu tiên, dữ liệu đầu vào được đưa vào hệ thống xử lý luồng, đôi khi được gọi là lớp thời gian thực. Lớp này chịu trách nhiệm chạy các tác vụ xử lý luồng và cung cấp khả năng xử lý dữ liệu thời gian thực. Sau đó, dữ liệu được chuyển đến lớp phục vụ (serving layer) để truy vấn bất kỳ kết quả cần thiết nào. Lưu ý rằng các thành phần của kiến trúc này không khác biệt đặc biệt so với kiến trúc Lambda. Hơn nữa, không cần phải trình bày thêm về bất kỳ công nghệ nào khác được sử dụng để triển khai kiến trúc này. Điều này là do kiến trúc Kappa có thể được triển khai bằng cách sử dụng cùng các công nghệ mã nguồn mở đã được trình bày trước đó để thực hiện kiến trúc Lambda, như được minh họa trong Hình 2.1.

Nhóm đã chọn triển khai hệ thống dự đoán giá tiền điện tử theo mô hình Kappa Architecture dựa trên các yếu tố kỹ thuật và yêu cầu nghiệp vụ đặc thù. Dữ liệu giao dịch tiền điện tử, với tốc độ cập nhật hàng nghìn transactions/giây trong thời điểm biến động mạnh, đòi hỏi xử lý streaming-first để đạt độ trễ dưới 1 giây, vì batch processing trong Lambda Architecture có thể gây trễ vài giờ, làm giảm độ chính xác dự đoán, như Makarov và Schoar et al đã phân tích [6]. (ii) Kiến trúc Kappa hỗ trợ reprocessing linh hoạt cho model retraining, error correction, và backtesting nhờ distributed commit log [8], đồng thời giảm 30-50% chi phí và effort bảo trì so với Lambda, theo Shahrivari et al [9]. (iii) Các công nghệ hiện đại như Kafka và Spark Structured Streaming cung

cấp exactly-once semantics, state management, và SQL interface, đáp ứng tốt ultra-low latency, 99.9% uptime, và scalability cho thị trường tiền điện tử 24/7. Do đó, Kappa Architecture không chỉ vượt trội về kỹ thuật mà còn phù hợp với yêu cầu xử lý real-time, cấu hình thấp, và độ trễ thấp của bài toán.

## 2.2. Apache Kafka - Xương sống của streaming platform

Apache Kafka là một nền tảng xử lý luồng dữ liệu phân tán, được thiết kế để xử lý hàng triệu sự kiện mỗi giây với độ trễ thấp. Trong hệ thống khuyến nghị theo kiến trúc Kappa, Kafka thường đóng vai trò là hệ thống trung gian chịu trách nhiệm thu thập và truyền dữ liệu người dùng theo thời gian thực từ front-end (website, app) đến các thành phần xử lý như Spark Streaming.



Hình 2.2. Kiến trúc cơ bản của Apache Kafka

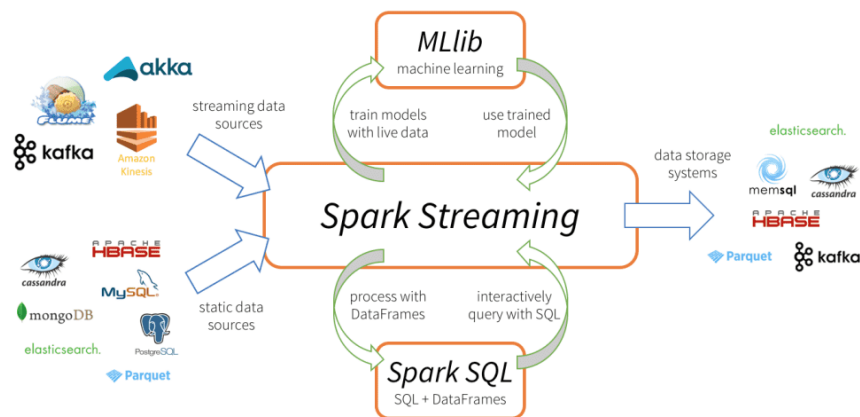
Nhờ khả năng mở rộng linh hoạt và độ tin cậy cao, Kafka giúp duy trì luồng dữ liệu ổn định giữa các thành phần, đảm bảo rằng mọi hành vi mới của người dùng đều được ghi nhận và phản ánh kịp thời trong hệ thống đề xuất. Cụ thể các ưu điểm như sau:

- **Khả năng chịu tải cao:** Có thể tăng cường khả năng xử lý bằng cách thêm các broker mới vào cụm xử lý luồng dữ liệu.
- **Đảm bảo tính nhất quán:** Sử dụng mô hình log-based để lưu các event, đảm bảo dữ liệu nguyên vẹn, luôn tuân theo thứ tự gửi.
- **Khả năng chịu lỗi cao:** Trong một cụm bao gồm N broker, cho phép tối đa N – 1 broker gặp sự cố mà vẫn đảm bảo dữ liệu được gửi đi.

- **Xử lý dòng sự kiện:** Hỗ trợ xử lý dòng sự kiện theo thời gian thực, giúp ứng dụng theo dõi và phản ứng nhanh chóng đối với các sự kiện quan trọng.

## 2.3. Apache Spark

Apache Spark là một nền tảng xử lý dữ liệu phân tán mã nguồn mở, được thiết kế để xử lý dữ liệu lớn với hiệu suất cao nhờ cơ chế tính toán trong bộ nhớ (in-memory computing). Spark hỗ trợ nhiều mô hình lập trình, bao gồm xử lý thời gian thực, học máy, và truy vấn SQL, đồng thời tích hợp dễ dàng với các công nghệ như Apache Kafka, AWS S3, Cassandra, và PyTorch, khiến nó trở thành lựa chọn lý tưởng cho các hệ thống dự đoán giá tiền điện tử theo thời gian thực.



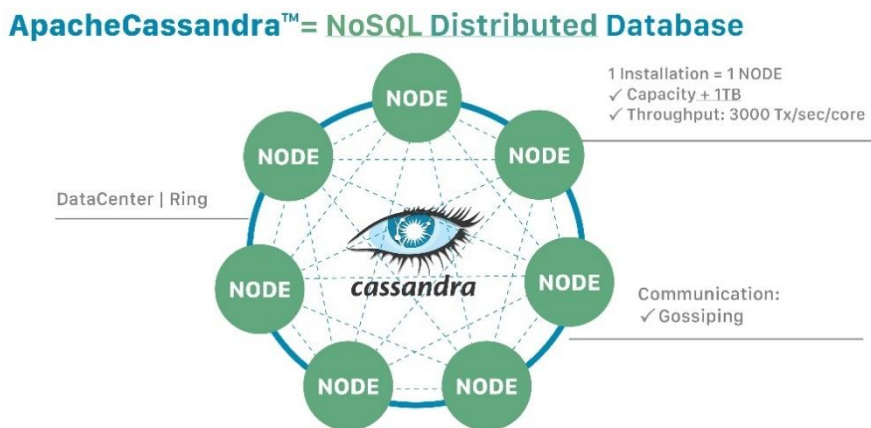
Hình 2.3. Apache Spark

Trong Kappa Architecture, Apache Spark đóng vai trò cốt lõi trong việc xử lý dữ liệu thời gian thực với Structured Streaming, một API mạnh mẽ cho phép xử lý luồng dữ liệu liên tục với độ trễ thấp và tính nhất quán cao. Dữ liệu giao dịch từ sàn Coinbase được truyền qua Kafka và xử lý bởi Structured Streaming để thực hiện các tác vụ như làm sạch dữ liệu, tính toán đặc trưng (features), và chuẩn bị đầu vào cho các mô hình học sâu. PySpark, thư viện Python của Spark, tích hợp với PyTorch thông qua TorchDistributor, cho phép huấn luyện phân tán các mô hình như LSTM, CNN-LSTM, và CNN-LSTM Attention trên cụm Spark [3]. TorchDistributor khởi tạo môi trường huấn luyện trên từng node, thiết lập giao tiếp giữa các worker để trao đổi gradient và cập nhật mô hình, tận dụng hạ tầng phân tán của Spark để tăng tốc độ huấn luyện. Structured Streaming cũng hỗ trợ lưu trữ dữ liệu đã xử lý vào Cassandra hoặc S3/MinIO, đảm bảo dữ liệu sẵn sàng cho

dự đoán giá liên tục (3 giờ tới) và trực quan hóa trên Grafana. Thiết kế streaming-native của Kappa Architecture, kết hợp với Spark, giúp hệ thống đáp ứng yêu cầu ultra-low latency (<1 giây) và high-throughput (10,000-100,000 messages/giây) của thị trường tiền điện tử.

## 2.4. Apache Cassandra

Apache Cassandra là một hệ quản trị cơ sở dữ liệu NoSQL phân tán mã nguồn mở, được thiết kế để xử lý khối lượng dữ liệu lớn với hiệu suất cao, khả năng mở rộng tuyến tính, và độ trễ thấp trên các hệ thống phân tán [1]. Với kiến trúc không có điểm lỗi đơn (no single point of failure) và hỗ trợ mô hình eventual consistency, Cassandra đặc biệt phù hợp cho các ứng dụng yêu cầu truy xuất dữ liệu thời gian thực, như dự đoán giá tiền điện tử.

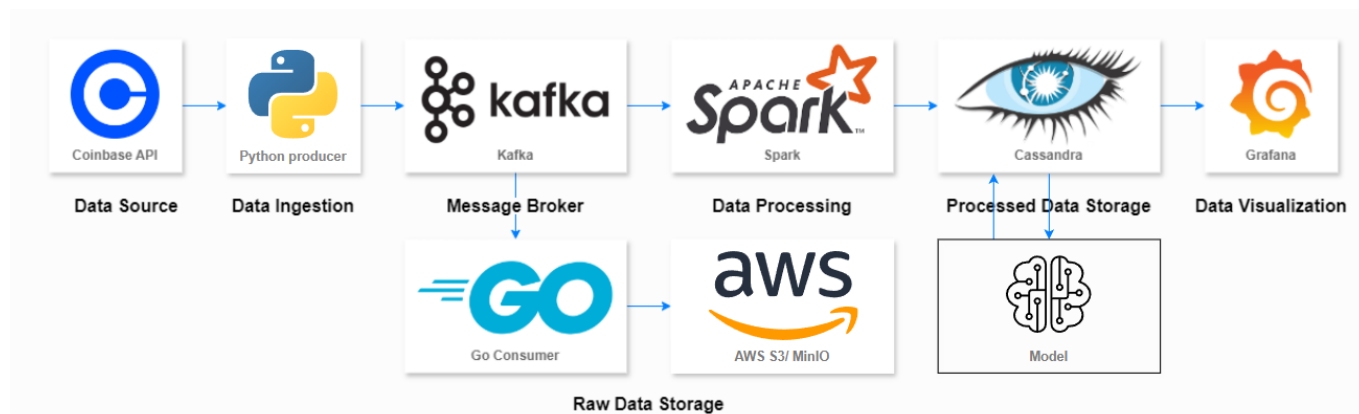


Hình 2.4. Apache Cassandra

Trong Kappa Architecture của hệ thống, Cassandra đóng vai trò lưu trữ dữ liệu giao dịch và lịch sử giá tiền điện tử từ luồng xử lý thời gian thực. Dữ liệu sau khi được xử lý bởi Structured Streaming được ghi vào Cassandra với độ trễ chỉ 10-30ms, đáp ứng yêu cầu ultra-low latency của thị trường tiền điện tử. Cassandra hỗ trợ truy vấn nhanh cho các tác vụ như cung cấp dữ liệu đầu vào cho mô hình dự đoán (LSTM, CNN-LSTM, CNN-LSTM Attention) và trực quan hóa xu hướng giá trên Grafana. Kiến trúc column-family của Cassandra cho phép lưu trữ hiệu quả dữ liệu chuỗi thời gian (time-series data), với timestamp chính xác đến microsecond, đáp ứng yêu cầu audit và compliance của thị trường tài chính. Ngoài ra, khả năng replication và multi-node đảm bảo tính chịu lỗi, duy trì 99.9% uptime ngay cả trong điều kiện biến động cao. Bằng cách tích hợp với Kafka và Structured Streaming, Cassandra tạo thành một phần không thể thiếu của pipeline dữ liệu, đảm bảo dữ liệu sẵn sàng cho dự đoán liên tục và phân tích thời gian thực.



## 2.5. Thiết kế pipeline dữ liệu



Hình 2.5. Pipeline dữ liệu thời gian thực

Hệ thống dự đoán giá tiền điện tử theo thời gian thực yêu cầu một pipeline dữ liệu mạnh mẽ để thu thập, truyền, lưu trữ, xử lý, dự đoán, và trực quan hóa dữ liệu giao dịch từ sàn Coinbase, đáp ứng độ trễ dưới 1 giây và thông lượng cao (10,000-100,000 messages/giây) trong thị trường biến động mạnh. Được thiết kế theo Kappa Architecture, pipeline sử dụng một luồng xử lý duy nhất để đảm bảo tính nhất quán, khả năng mở rộng, và hiệu suất tối ưu, tránh sự phức tạp của kiến trúc Lambda.

Việc thiết kế pipeline dựa trên ba nguyên tắc cốt lõi. Thứ nhất, nguyên tắc tách biệt lỏng lẻo (loose coupling) giữa các thành phần, cho phép mỗi service có thể phát triển và mở rộng độc lập. Thứ hai, khả năng chịu lỗi cao thông qua việc nhân bản dữ liệu và cơ chế phục hồi tự động. Thứ ba, khả năng mở rộng linh hoạt để đáp ứng với sự biến động của thị trường tiền điện tử, nơi khối lượng giao dịch có thể tăng đột biến trong thời gian ngắn.

Các thành phần chính bao gồm Coinbase WebSocket API, Java Kafka Producer, Kafka Broker, Go Kafka Consumer, AWS S3/MinIO, Spark Structured Streaming, Cassandra, Grafana, và dịch vụ dự đoán, được triển khai thông qua Docker Compose. Pipeline phân chia dữ liệu thành hai luồng: một luồng trực quan hóa giá thời gian thực qua Grafana và một luồng dự đoán giá liên tục (3 giờ tới) bằng các mô hình LSTM, CNN-LSTM, CNN-LSTM Attention sẽ được giới thiệu ở dưới.

### 2.5.1. Thu thập dữ liệu - Coinbase WebSocket API

Lớp thu thập dữ liệu sử dụng Coinbase WebSocket API để cung cấp luồng dữ liệu thị trường theo thời gian thực với độ trễ cực thấp. Giao thức WebSocket được chọn thay vì API REST truyền thống nhờ khả năng duy trì kết nối liên tục và mô hình giao tiếp dựa trên cơ chế đẩy (push-based), giúp giảm chi phí của các lần bắt tay HTTP lặp lại.

Coinbase Advanced Trade WebSocket API cung cấp nhiều kênh dữ liệu khác nhau, phục vụ các loại dữ liệu riêng biệt:

1. Kênh Ticker: Cập nhật giá theo thời gian thực, bao gồm thông tin về khối lượng giao dịch và biến động giá.
2. Kênh Trades: Ghi nhận các giao dịch riêng lẻ với dấu thời gian tính bằng micro giây.
3. Kênh Level2: Cập nhật sổ lệnh (order book) để phân tích độ sâu thị trường.
4. Kênh Candles: Dữ liệu OHLCV (open, high, low, close, volume) theo nhiều khung thời gian.

Trong triển khai hiện tại, hệ thống đăng ký (subscribe) vào các kênh Ticker và Candles cho ba cặp tiền điện tử chính: BTC-USD, ETH-USD, và XRP-USD. Định dạng thông điệp tuân theo lược đồ JSON với thứ tự đảm bảo theo product\_id, điều này rất quan trọng để duy trì tính nhất quán của dữ liệu trong quá trình xử lý tiếp theo.

### 2.5.2. Kafka Producer Service – Hàm truyền dữ liệu tới Kafka

Dịch vụ Kafka Producer được triển khai như một microservice độc lập, đóng vai trò cầu nối giữa nguồn dữ liệu WebSocket và hệ thống truyền dữ liệu Kafka. Thiết kế của dịch vụ tập trung vào ba yếu tố chính: **độ tin cậy**, **hiệu năng**, và **khả năng giám sát**.

1. **Độ tin cậy**: Producer tích hợp cơ chế quản lý kết nối với khả năng tự động kết nối lại khi xảy ra sự cố. Chiến lược **exponential backoff** được áp dụng để tránh gây áp lực lên máy chủ trong trường hợp mất kết nối kéo dài. Mỗi thông điệp nhận được đều được kiểm tra lược đồ (schema validation) nghiêm ngặt trước khi gửi, đảm bảo chất lượng dữ liệu và ngăn chặn dữ liệu lỗi lan truyền trong hệ thống.

2. **Hiệu năng:** Producer sử dụng kỹ thuật **batching** (gộp nhóm) và **nén dữ liệu** để tối ưu hóa thông lượng. Các thông điệp được nhóm lại theo thời gian hoặc kích thước trước khi gửi, giảm số lượng lệnh gọi mạng. Thuật toán nén **Snappy** được chọn vì cân bằng tốt giữa tỷ lệ nén và tốc độ xử lý, phù hợp với yêu cầu thời gian thực của hệ thống.
3. **Cách hoạt động:** Dữ liệu từ Coinbase được truyền đến Kafka Broker thông qua Kafka Producer, được triển khai trong tệp **producer.py** bằng thư viện **kafka-python**. Producer gửi dữ liệu **Ticker** vào topic **coin-data** và dữ liệu **Candles** vào topic **coin-data-model**, với khóa phân vùng là **product\_id** (ví dụ: BTC-USD). Kafka, một nền tảng truyền dữ liệu phân tán, lưu trữ dữ liệu trong các topic được phân vùng và nhân bản, hỗ trợ thông lượng cao (10.000-100.000 thông điệp/giây trong giờ cao điểm) và khả năng chịu lỗi. Tệp **docker-compose.yml** cấu hình Kafka và **kafka-init** để tạo hai topic với 3 phân vùng và hệ số nhân bản (replication factor) là 1. Producer mã hóa thông điệp để đảm bảo tính nhất quán và giảm lỗi phân tích cú pháp. Kafka cho phép xử lý lại dữ liệu để huấn luyện mô hình hoặc sửa lỗi, một ưu điểm chính của kiến trúc Kappa.

### 2.5.3. Apache Kafka - xương sống của streaming

Apache Kafka đóng vai trò trung tâm trong việc phân phối dữ liệu với khả năng xử lý thông lượng cao và đảm bảo độ tin cậy. Cấu hình Kafka cluster được tối ưu cụ thể cho workload của cryptocurrency data với nhiều cân nhắc kỹ lưỡng.

**Thiết kế topic:** Các topic trong Kafka được thiết kế để tách biệt logic giữa các loại dữ liệu:

1. Topic coin-data chứa dữ liệu Ticker tần suất cao với thời gian lưu trữ 7 ngày, đủ để phục vụ phân tích ngắn hạn và cung cấp dữ liệu thời gian thực cho bảng điều khiển giám sát.
2. Topic coin-data-model lưu trữ dữ liệu Candles đã được tổng hợp theo các khung thời gian, phục vụ trực tiếp cho việc huấn luyện và suy luận của mô hình dự đoán.
3. Kết quả dự đoán không được truyền qua Kafka mà được ghi trực tiếp vào Cassandra để đảm bảo truy cập với độ trễ thấp cho lớp trực quan hóa.

**Chiến lược phân vùng:** Mỗi topic được cấu hình với 3 phân vùng, con số này được xác định dựa trên công thức cân bằng giữa thông lượng mong đợi và khả năng xử lý của consumer.

Phân vùng dựa trên khóa (key-based partitioning) đảm bảo các thông điệp của cùng một product\_id luôn được định tuyến đến cùng một phân vùng, duy trì thứ tự xử lý cho mỗi cặp giao dịch.

**Cấu hình nhân bản:** Hệ số nhân bản là 3 và số lượng bản sao đồng bộ tối thiểu là 2, tạo ra sự cân bằng giữa độ bền dữ liệu và hiệu suất ghi. Cấu hình này cho phép hệ thống chịu được tối đa 2 lỗi broker mà không mất dữ liệu, đồng thời duy trì khả năng ghi khi một broker gặp sự cố.

Sau khi nhận được data, Hệ thống triển khai mô hình dual consumer để phục vụ hai nhu cầu xử lý khác nhau: lưu trữ dữ liệu thô và xử lý thời gian thực cho phân tích.

#### 2.5.4. Lưu trữ dữ liệu thô

Go Consumer cho lưu trữ dữ liệu thô: Được thiết kế để truyền dữ liệu từ Kafka sang kho lưu trữ đối tượng (object storage) với thông lượng tối đa. Ngôn ngữ Go được chọn nhờ khả năng xử lý đồng thời vượt trội với goroutines và channels, cho phép xử lý hàng nghìn thông điệp cùng lúc với dấu chân bộ nhớ thấp. Consumer áp dụng mô hình worker pool với số lượng worker có thể điều chỉnh theo khối lượng công việc. Các thông điệp được gộp nhóm theo cửa sổ thời gian hoặc ngưỡng kích thước trước khi tải lên S3/MinIO, tối ưu hóa số lượng lệnh gọi API và giảm chi phí lưu trữ.

Go Kafka Consumer được triển khai trong tệp consumer.go bằng ngôn ngữ Go để tối ưu hiệu suất, đọc dữ liệu từ cả hai topic coin-data và coin-data-model, sau đó ghi vào AWS S3/MinIO. Tệp consumer.go sử dụng AWS SDK để kết nối với MinIO (cấu hình qua AWS\_ENDPOINT, AWS\_ACCESS\_KEY\_ID, AWS\_S3\_FORCE\_PATH\_STYLE), lưu dữ liệu dưới dạng tệp JSON theo cấu trúc data\_type/product\_id/timestamp.json (ví dụ: ticker/BTC-USD/1623456789.json). Tệp docker-compose.yml định nghĩa dịch vụ minio với bucket ie212-coinbase-data được tạo tự động qua minio-setup. Độ trễ ghi vào S3/MinIO dao động từ 10-30ms, đảm bảo lưu trữ nhanh chóng và hỗ trợ các tác vụ như kiểm toán, kiểm tra ngược (backtesting), hoặc khôi phục dữ liệu. MinIO cung cấp độ bền cao và khả năng mở rộng tuyến tính, đáp ứng yêu cầu lưu trữ dài hạn của thị trường tài chính.

### 2.5.5. Xử lý dữ liệu thời gian thực

Spark Structured Streaming Consumer thực hiện các tác vụ xử lý phức tạp với nhiều bước biến đổi và tổng hợp dữ liệu. Mô hình xử lý micro-batch được cấu hình với khoảng thời gian kích hoạt 100ms, cung cấp khả năng xử lý gần thời gian thực với hiệu quả tài nguyên tối ưu. Quản lý trạng thái được xử lý tự động bởi Spark thông qua cơ chế checkpointing, đảm bảo ngữ nghĩa xử lý đúng một lần (exactly once) khi kết hợp với hỗ trợ giao dịch của Kafka.

#### Luồng dữ liệu:

1. Dữ liệu từ Kafka được xử lý bởi Spark Structured Streaming, triển khai trong tệp `spark_processor.py` bằng PySpark.
2. Structured Streaming đọc dữ liệu từ các topic `coin-data` và `coin-data-model` thông qua Kafka connector, áp dụng các phép biến đổi thời gian thực bằng mô hình micro-batch với độ trễ 50-200ms.
3. Phân tích JSON: Dữ liệu Kafka được chuyển đổi từ chuỗi JSON thành DataFrame bằng hàm `from_json`, sau đó được làm sạch và chuẩn hóa.
4. Dữ liệu đã xử lý được ghi vào Cassandra (các bảng `prices` và `candles`) với cơ chế checkpointing để đảm bảo khả năng phục hồi.

### 2.5.6. Apache Cassandra - Lưu trữ và truy xuất dữ liệu

Cassandra được tối ưu hóa cho khối lượng công việc chuỗi thời gian của tiền điện tử. Thiết kế lược đồ cơ sở dữ liệu tuân theo các phương pháp tốt nhất cho dữ liệu ghi nhiều, sắp xếp theo thời gian:

1. **Bảng prices:**
  - Khóa phân vùng (partition key): `product_id`, đảm bảo phân phối đều trên các node.
  - Cột sắp xếp (clustering column): `time`, sắp xếp tự nhiên theo thời gian.
  - Sắp xếp DESC: Dữ liệu gần nhất được ưu tiên (phù hợp với mẫu truy vấn phổ biến).
2. **Bảng candles:** Được thiết kế đặc biệt để khớp với dữ liệu OHLCV đã được huấn luyện trong mô hình.
3. **Bảng predictions:** Sử dụng khóa phân vùng tổng hợp để đảm bảo phân phối đều và hiệu quả truy vấn cho các dự đoán dành riêng cho mô hình, cập nhật 5 phút/lần.

Về luồng dữ liệu, Apache Cassandra lưu trữ dữ liệu đã xử lý và kết quả dự đoán trong keyspace coinbase, với các bảng prices, candles và prediction. File docker-compose.yml định nghĩa dịch vụ cassandra với image tùy chỉnh (cassandra.Dockerfile) và cấu hình keyspace coinbase. Cassandra sử dụng mô hình column-family, tối ưu cho dữ liệu chuỗi thời gian với độ trễ ghi/đọc 10-30ms. Dịch vụ dự đoán (price-predictor) và Grafana truy vấn dữ liệu từ Cassandra với độ trễ dưới 100ms, hỗ trợ dự đoán liên tục và trực quan hóa. Replication và multi-node đảm bảo uptime 99.9% và khả năng chịu lỗi.

### 2.5.7. Dự đoán và trực quan hóa

Grafana cung cấp khả năng trực quan hóa toàn diện với nhiều bảng điều khiển phục vụ các nhu cầu khác nhau của người dùng.

#### Kiến trúc bảng điều khiển:

1. Bảng điều khiển giá thời gian thực (dashboard.json):

- Biểu đồ giá trực tiếp cho tất cả các cặp giao dịch.
- Tần suất làm mới 5 giây để mang lại cảm giác thời gian thực.
- Các bảng chuỗi thời gian với trục Y tự động điều chỉnh.



Hình 2.6. Visualize giá tiền điện tử theo thời gian thực

2. Bảng điều khiển dự đoán (predict\_dashboard.json):

- Hiển thị kép: Giá thực tế so với giá dự đoán.
- Hiển thị nhiều khung thời gian.
- Trực quan hóa khoảng tin cậy.
- Các bảng hiển thị chỉ số hiệu suất mô hình.



Hình 2.7. Visualize giá tiền điện tử dự đoán theo thời gian thực

Dịch vụ dự đoán sử dụng mô hình LSTM để dự đoán giá tiền điện tử trong 3 giờ tới dựa trên dữ liệu thời gian thực từ Coinbase, được tích hợp vào pipeline theo Kappa Architecture để đảm bảo độ trễ thấp. Batch prediction xử lý dữ liệu theo lô định kỳ (vài phút đến giờ), phù hợp cho phân tích dài hạn nhưng không đáp ứng thị trường tiền điện tử biến động nhanh. Real-time inference, được chọn cho hệ thống, dự đoán ngay khi dữ liệu đến, sử dụng dữ liệu OHLC từ Cassandra để hỗ trợ quyết định giao dịch tức thì. price-predictor chạy inference mỗi 5 phút (`PREDICTION_INTERVAL=300`), đáp ứng yêu cầu soft real-time

Về luồng dữ liệu dự đoán:

1. Dịch vụ **price-predictor**, được triển khai trong **docker-compose.yml**, truy xuất dữ liệu Candles từ bảng candles trong Cassandra để cung cấp đầu vào cho các mô hình LSTM, CNN-LSTM, CNN-LSTM Attention (tệp mô hình tại `/app/model/checkpoints/best_epoch_16.pt`).
2. Mô hình dự đoán giá trong 3 giờ tới với độ trễ 100-500ms, sử dụng cấu hình từ `train_config.yaml` và chạy trên CPU.
3. Kết quả dự đoán được ghi vào bảng predictions trong Cassandra, với các trường như `product_id`, `model_name` (ví dụ: `lstm_model_v1`), và `predicted_price`.

4. Grafana truy vấn các bảng prices, candles, và predictions để hiển thị trên bảng điều khiển. Cả hai luồng đảm bảo độ trễ end-to-end dưới 1 giây, phù hợp với đặc thù thị trường tiền điện tử.

#### 2.5.8. Triển khai hệ thống:

Docker Compose điều phối toàn bộ hệ thống cho quá trình phát triển và kiểm thử:

- Phụ thuộc dịch vụ: Các khai báo depends\_on đảm bảo thứ tự khởi động đúng. Kiểm tra sức khỏe (health checks) xác nhận dịch vụ sẵn sàng trước khi các dịch vụ phụ thuộc khởi động.
- Quản lý tài nguyên: Giới hạn bộ nhớ và CPU ngăn chặn tình trạng cạn kiệt tài nguyên. Chính sách khởi động lại đảm bảo tự động khôi phục từ sự cố.
- Mạng: Mạng cầu tùy chỉnh cho phép khám phá dịch vụ qua DNS. Ánh xạ cổng cho phép truy cập bên ngoài khi cần.

Docker Compose đơn giản hóa việc thiết lập và kiểm thử trên máy đơn, với các volume để lưu trữ dữ liệu bền vững. Cấu hình có thể chuyển sang Kubernetes (GKE, EKS) bằng cách bổ sung pod scaling, replica sets, và load balancer, hỗ trợ triển khai sản xuất với các đợt tăng lưu lượng đột biến.

Tóm lại, luồng dữ liệu bắt đầu bằng việc thu thập sự kiện thị trường thô từ Coinbase qua WebSocket API, sau đó Producer xác thực, chuẩn hóa và gửi dữ liệu đến các topic Kafka tương ứng. Kafka phân phối thông điệp qua các phân vùng, với các consumer phối hợp xử lý và quản lý offset để đảm bảo không mất dữ liệu. Dữ liệu thô lưu trữ trên S3/MinIO thông qua consumer go, còn consumer Spark Structured Streaming đọc dữ liệu theo micro-batch, thực hiện biến đổi, quản lý trạng thái và xử lý dữ liệu trễ bằng watermark. Dữ liệu đã xử lý được ghi vào Cassandra, đồng thời mô hình dự đoán truy xuất dữ liệu từ Cassandra để dự đoán giá liên tục, với kết quả được lưu lại và tối ưu hóa chỉ mục để truy vấn nhanh. Cuối cùng, Grafana truy vấn dữ liệu từ Cassandra, lưu kết quả vào bộ nhớ đệm và hiển thị trực quan giá thực tế cùng giá dự đoán trên bảng điều khiển với cơ chế tự động làm mới để duy trì tính thời gian thực.

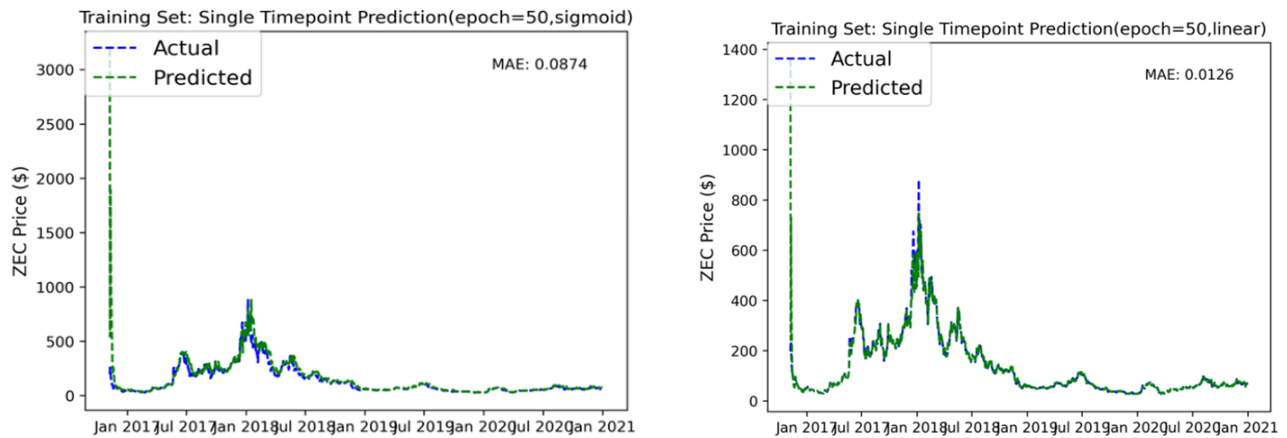


## Chương 3. CÁC CÔNG VIỆC LIÊN QUAN

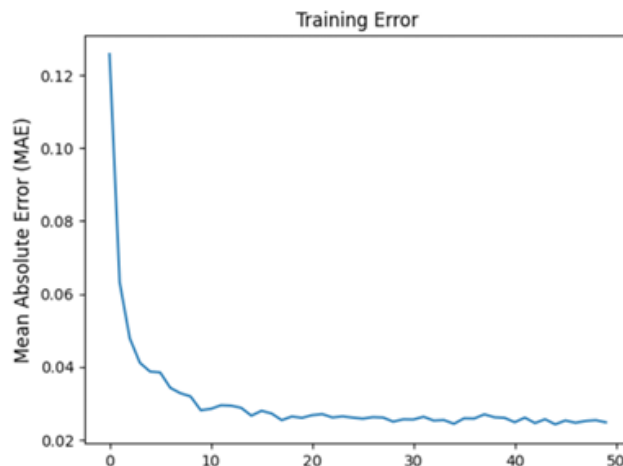
### 3.1. Các phương pháp dự đoán giá tiền điện tử

#### Paper: Comparative Analysis of LSTM, GRU and Transformer Deep Learning Models for Cryptocurrency ZEC Price Prediction Performance

Tác giả tập trung vào dự báo giá đồng Zcash (ZEC) bằng các mô hình học sâu phổ biến. Mục tiêu là so sánh khả năng dự báo của ba kiến trúc mạng nơ-ron: LSTM, GRU và Transformer trên dữ liệu giá tiền mã hóa. Nghiên cứu giới thiệu hai đặc trưng kỹ thuật mới (“close\_off\_high” và “volatility”) và phân tích mối tương quan của chúng với giá ZEC. Kết quả cho thấy LSTM và GRU vượt trội hơn hẳn so với Transformer trong việc dự báo xu hướng giá ZEC. Tác giả cũng khảo sát ảnh hưởng của hàm kích hoạt (kết luận rằng hàm tuyến tính tốt hơn sigmoid) và thảo luận các thách thức như quá khớp, đa cộng tuyến, đồng thời chỉ ra hạn chế của việc chỉ nghiên cứu trên một loại tiền mã hóa.



Hình 3.1. Đánh giá training test



Hình 3.2. Đánh giá training error

Nghiên cứu trên có nhiều điểm có thể tham khảo cho đề tài dự báo giá tiền mã hóa thời gian thực (BTC/USD) từ Coinbase của nhóm. Cụ thể:

**Đặc trưng dữ liệu:** Ý tưởng tạo hai đặc trưng “close\_off\_high” và “volatility” đã được chứng minh mang lại thông tin hữu ích cho mô hình. Hệ thống của bạn có thể áp dụng tương tự cho dữ liệu BTC/USD để phản ánh các xu hướng giá và độ biến động hàng ngày. Việc kiểm tra hệ số tương quan giữa các đặc trưng này và giá BTC cũng sẽ giúp hiểu quan hệ dữ liệu đầu vào-chuỗi mục tiêu.

**Mô hình tham khảo:** Kết quả cho thấy GRU hoạt động rất tốt (MAE thấp nhất) trong dự báo giá ZEC. Do đó, ngoài các mô hình LSTM hiện có, bạn có thể mở rộng thêm mô hình GRU (hoặc kết hợp CNN+GRU) vào hệ thống của mình. Việc này có thể cải thiện độ chính xác dự báo, bởi GRU thường huấn luyện nhanh và tránh quá khớp tốt hơn LSTM trong nhiều trường hợp.

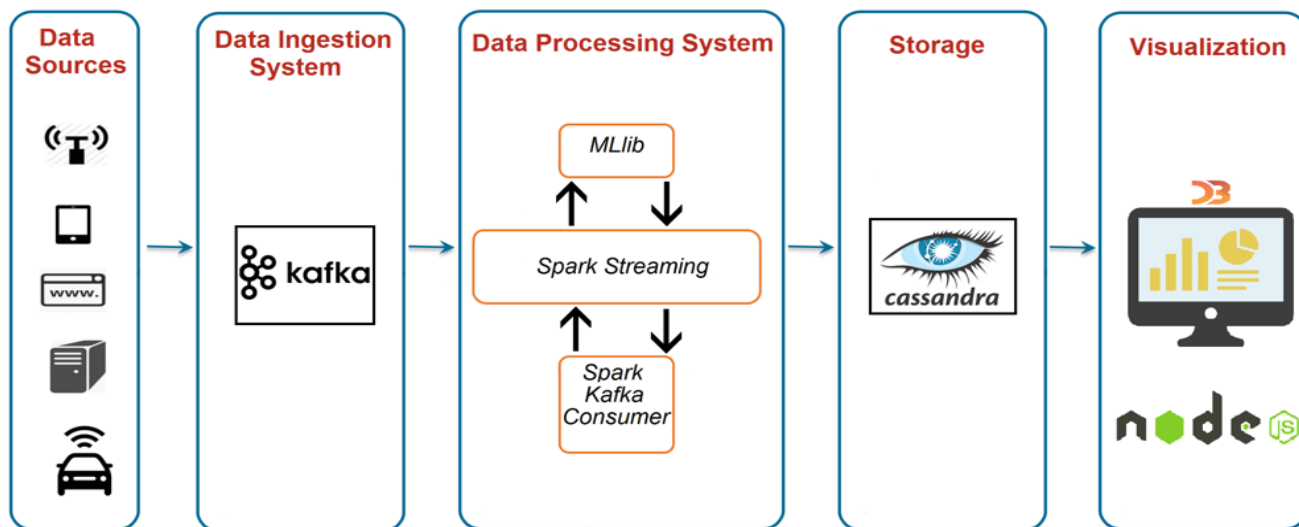
**Kiến trúc chú ý:** Dù nghiên cứu này cho thấy Transformer đơn giản cho kết quả thấp hơn, cơ chế tự chú ý vẫn có tiềm năng. Hệ thống của bạn đã dùng CNN-LSTM-Attention, đây là cách kết hợp sức mạnh của CNN và attention tương tự ý tưởng của Transformer. Bạn có thể nghiên cứu tiếp tục một mô hình Transformer thuần túy (hoặc các biến thể như Informer, Autoformer dành riêng cho dữ liệu thời gian) để so sánh hiệu quả. Ngay cả nghiên cứu này cũng khuyến nghị không bỏ qua Transformer hoàn toàn mà cần điều chỉnh tham số và kiến trúc cho phù hợp hơn.

**Công thức huấn luyện:** Các thiết lập về hàm mất mát MAE, optimizer Adam, dropout, batch size,... trong nghiên cứu này là gợi ý tham số ban đầu để thử nghiệm. Ví dụ, hệ thống của bạn có thể áp dụng MAE và thử nghiệm với hàm kích hoạt đầu ra tuyến tính như gợi ý trong bài. Đồng thời, lưu ý sử dụng kỹ thuật regularization và early stopping để tránh quá khớp, nhất là với dữ liệu streaming có thể chứa nhiễu.

**Mở rộng nghiên cứu:** Kết quả khuyến nghị rằng nên xem xét nhiều tiền mã hóa hơn và dữ liệu bổ trợ. Trong tương lai, hệ thống của bạn có thể tích hợp thêm dữ liệu on-chain (thanh khoản, giao dịch), dữ liệu cảm tính thị trường (tin tức, mạng xã hội) như đề xuất của tác giả. Điều này giúp hệ thống dự báo BTC/USD có bối cảnh rộng hơn, cải thiện độ chính xác và độ tin cậy.

### 3.2. Pipeline dữ liệu trong dự đoán tài chính

Trong bài báo “Real-time Text Analytics Pipeline Using Open-source Big Data Tools” (Nazeer et al., 2017), các tác giả đề xuất một kiến trúc phân tích văn bản thời gian thực sử dụng hoàn toàn các công cụ mã nguồn mở gồm **Apache Kafka**, **Apache Spark Streaming**, **Apache Cassandra** và thư viện **D3.js**. Mục tiêu chính là giảm độ trễ khi xử lý dòng dữ liệu văn bản lớn, ví dụ như các dòng tweet từ Twitter, nhằm phục vụ phân tích cảm xúc gần như tức thời.



Hình 3.3. Cấu trúc cơ bản của bài báo Nazeer et al., 2017

Hình 3.2: Minh họa kiến trúc pipeline phân tích văn bản thời gian thực đề xuất trong nghiên cứu.

Dữ liệu văn bản (ở đây là các tweet) được thu thập từ nhiều nguồn và đưa vào **Apache Kafka** làm hệ thống nhắn tin phân tán (data ingestion). Kafka chịu trách nhiệm tiếp nhận luồng dữ liệu đầu vào và phát cho các thành phần tiếp theo. **Apache Spark Streaming** nhận dữ liệu từ Kafka và xử lý trực tuyến; tác giả sử dụng MLlib của Spark để áp dụng mô hình học máy trên luồng dữ liệu, cụ thể là phân tích cảm xúc với thuật toán Naive Bayes. Kết quả xử lý (ví dụ nhãn cảm xúc của mỗi tweet) được lưu trữ vào **Apache Cassandra** – một cơ sở dữ liệu NoSQL phân tán, đáp ứng yêu cầu ghi và đọc hiệu năng cao cho dữ liệu lớn. Cuối cùng, kết quả phân tích được đưa đến thành phần trực quan hóa: bài báo triển khai ứng dụng web Node.js với thư viện D3 để hiển thị bảng điều khiển (dashboard) trực quan các thống kê theo thời gian thực. Như vậy, pipeline đề xuất bao gồm chuỗi: nguồn dữ liệu → Kafka (ingestion) → Spark Streaming (xử lý – xử lý cảm xúc với MLlib) → Cassandra (lưu trữ) → Dashboard (Node/D3), đảm bảo luồng dữ liệu liên tục và độ trễ thấp.

- **Apache Kafka:** Hệ thống nhắn tin phân tán theo mô hình publish-subscribe, chịu trách nhiệm ingest dữ liệu tweet với tốc độ cao.

- **Apache Spark Streaming:** Xử lý luồng dữ liệu ngay tại bộ nhớ (in-memory processing). Spark lắng nghe topic từ Kafka (Spark Kafka Consumer), áp dụng mô hình học máy (MLlib) để phân loại cảm xúc tweet thành tích cực/tiêu cực.
- **Apache Cassandra:** Cơ sở dữ liệu NoSQL phi tập trung, lưu trữ kết quả xử lý (nhãn cảm xúc và các thông kê liên quan). Cassandra được dùng vì khả năng mở rộng linh hoạt, chịu lỗi tốt và độ trễ thấp khi ghi/đọc.
- **Trực quan hóa (Node.js + D3):** Kết quả từ Cassandra được truy xuất bởi ứng dụng web viết bằng Node.js, dùng thư viện D3 để vẽ biểu đồ hiển thị số lượng từ khóa phổ biến trong các tweet tích cực/tiêu cực theo thời gian. Giao diện cho phép người dùng tìm kiếm từ khóa và so sánh tần suất xuất hiện giữa các nhóm tweet.

Hệ thống của nhóm có nhiều điểm tương đồng với pipeline được trình bày trong nghiên cứu của Nazeer et al., đặc biệt về mặt kiến trúc và công nghệ lõi: đều sử dụng **Apache Kafka** để tiếp nhận dữ liệu thời gian thực, **Apache Spark** để xử lý và tính toán, và **Apache Cassandra** làm lớp lưu trữ phân tán. Thành phần trực quan hóa cũng hiện diện trong cả hai hệ thống – dù sử dụng công nghệ khác nhau (Grafana trong hệ thống của tôi so với Node.js/D3 trong nghiên cứu), nhưng đều nhằm mục đích cung cấp dashboard phân tích theo thời gian thực.

Mặc dù khác biệt về mục tiêu – nghiên cứu gốc hướng đến phân tích cảm xúc từ dữ liệu văn bản, còn hệ thống của nhóm tập trung vào **dự báo giá BTC/USD** – song cách tiếp cận kiến trúc có thể áp dụng song song. Nghiên cứu của Nazeer et al. cung cấp **cơ sở thực nghiệm đáng tin cậy**, trong đó việc triển khai trên các cụm có quy mô khác nhau (1–3 VM) và đo lường các chỉ số như độ trễ, throughput là gợi ý thiết thực cho việc đánh giá hiệu năng hệ thống của nhóm.

Cụ thể, nhóm có thể áp dụng phương án mở rộng cụm Spark và Cassandra giúp giảm độ trễ xử lý, gợi ý hướng tối ưu hiệu quả khi quy mô dữ liệu (như khối lượng giao dịch hoặc độ sâu thị trường) tăng lên.

Hơn nữa, bài báo cũng nhấn mạnh đến **auto-scaling** – khả năng mở rộng linh hoạt hệ thống theo mức tải, đây là một tính năng quan trọng có thể tích hợp vào hệ thống của nhóm trong các giai đoạn cải tiến tiếp theo.

### 3.3. Khoảng trống nghiên cứu

Các nghiên cứu trước đây về dự đoán giá tiền điện tử thường gặp một số hạn chế, làm giảm hiệu quả trong bối cảnh thị trường biến động nhanh. Nhiều công trình, như nghiên cứu của Li et al. (2020), tập trung vào dự đoán giá bằng mô hình học sâu (LSTM, GRU) nhưng chủ yếu sử dụng dữ liệu lịch sử và xử lý theo lô (batch processing), thiếu khả năng dự đoán thời gian thực cần thiết cho giao dịch tức thì [1].

Một số nghiên cứu khác, như Nazeer et al. (2017), phát triển các hệ thống trực quan hóa dữ liệu thị trường nhưng không tích hợp dự đoán giá, hạn chế khả năng hỗ trợ quyết định [2]. Ngoài ra, các hệ thống thường không khai thác dữ liệu từ Coinbase, một sàn giao dịch lớn với luồng dữ liệu đáng tin cậy, hoặc không tận dụng giao diện trực quan hóa tương tác như Grafana để cung cấp thông tin theo thời gian thực [3].

Đề án của chúng ta giải quyết các khoảng trống này thông qua:

- **Pipeline thời gian thực:** Sử dụng Coinbase WebSocket API, Kafka, và Spark Structured Streaming để thu thập, truyền, và xử lý dữ liệu với độ trễ dưới 1 giây, đáp ứng yêu cầu soft real-time của thị trường tiền điện tử [4].
- **Mô hình học sâu tối ưu:** Tích hợp các mô hình LSTM, CNN-LSTM, và CNN-LSTM Attention trong dịch vụ price-predictor, dự đoán giá liên tục (3 giờ tới) với độ trễ 100-500ms, vượt trội so với batch prediction [5].
- **Giao diện Grafana:** Kết hợp dự đoán và trực quan hóa trên các dashboard tương tác, hiển thị giá thực tế và giá dự đoán (bảng predictions trong Cassandra), cải thiện khả năng theo dõi và ra quyết định.

Bằng cách kết hợp dự đoán thời gian thực và trực quan hóa tương tác, đề án không chỉ khắc phục hạn chế của các nghiên cứu trước mà còn cung cấp một giải pháp toàn diện, hỗ trợ hiệu quả cho giao dịch và phân tích thị trường tiền điện tử.

## Chương 4. PHƯƠNG PHÁP ÁP DỤNG CHO BÀI TOÁN

Nội dung nghiên cứu và phương pháp luận đóng vai trò trung tâm của đề tài, chú trọng vào việc trình bày chi tiết về cách mà nghiên cứu được thực hiện và phương pháp được sử dụng để đạt được mục tiêu đề ra. Chương này không chỉ làm rõ chi tiết về dữ liệu và các tài nguyên phương pháp tiền xử lý dữ liệu và lý thuyết của mô hình thực nghiệm, mà còn giúp độc giả hiểu rõ cách tiếp cận vấn đề nghiên cứu.

### 4.1. Mô tả dữ liệu

Dữ liệu lịch sử giá Bitcoin (BTC) theo cặp giao dịch BTC/USDT đã được thu thập từ nền tảng Binance trong khoảng thời gian từ ngày **17 tháng 08 năm 2017** đến ngày **5 tháng 05 năm 2025**. Dữ liệu được ghi nhận theo từng khung thời gian **5 phút**, phản ánh đầy đủ và chi tiết sự biến động giá của Bitcoin trong suốt hơn bảy năm qua. Với tần suất ghi nhận dày đặc, mỗi ngày có tới

đa 288 bản ghi (tương ứng với 288 phiên giao dịch 5 phút), tạo thành một tập dữ liệu lớn và giàu thông tin, phù hợp để phục vụ cho việc phân tích và dự đoán.

<b>Thời gian giao dịch</b>	<b>Giá mở cửa</b>	<b>Giá cao nhất</b>	<b>Giá thấp nhất</b>	<b>Giá đóng cửa</b>	<b>Khối lượng giao dịch</b>
8/17/2017 4:00:00 AM	4261.48	4280.56	4261.48	4261.48	2.189061
8/17/2017 4:05:00 AM	4261.48	4261.48	4261.48	4261.48	0
...	...	...	...	...	...
12/31/2024 11:50:00 PM	187.93	189.90	189.25	189.25	51,216,800
12/31/2024 11:55:00 PM	93646.9 7	93676.98	93576	93576	17.03553

Bảng 4.1 Bộ dữ liệu giá BTC-USD

Bộ dữ liệu bao gồm các trường thông tin chính sau:

- **Thời gian giao dịch:** Ghi nhận thời điểm bắt đầu và kết thúc của mỗi phiên giao dịch 5 phút.
- **Giá mở cửa:** Mức giá đầu tiên được ghi nhận tại thời điểm bắt đầu phiên.
- **Giá cao nhất:** Mức giá cao nhất được ghi nhận trong suốt phiên 5 phút.
- **Giá thấp nhất:** Mức giá thấp nhất được ghi nhận trong cùng phiên.
- **Giá đóng cửa:** Mức giá cuối cùng trước khi kết thúc phiên 5 phút, phản ánh giá trị gần nhất tại thời điểm đó.
- **Khối lượng giao dịch:** Số lượng Bitcoin được giao dịch trong phiên, thể hiện mức độ quan tâm và hoạt động của thị trường.

Trong lĩnh vực giao dịch tiền điện tử, mức giá được ghi nhận sau cùng trong một phiên là chỉ số thể hiện rõ nhất giá trị thị trường tại thời điểm đó. Đặc biệt với các khung thời gian ngắn như 5

phút, giá đóng cửa đóng vai trò như một tín hiệu quan trọng để nhà đầu tư theo dõi biến động ngắn hạn và xác định điểm vào hoặc thoát lệnh hiệu quả.

Để cải thiện hiệu quả dự báo, đặc biệt là trong việc nhận diện xu hướng ngắn hạn và biến động thị trường, nhiều đặc trưng đã được trích xuất từ dữ liệu gốc. Những đặc trưng này bao gồm::

- **Tỷ suất sinh lợi theo logarit (log return):** Được tính để phản ánh sự thay đổi tương đối của giá qua từng khoảng thời gian, mang lại cái nhìn rõ nét về động thái giá cả. Việc sử dụng log return giúp làm mượt chuỗi thời gian, giảm thiểu tác động của các biến động cực đại và tạo ra phân phối gần chuẩn, điều này đặc biệt phù hợp cho các mô hình học máy vốn yêu cầu dữ liệu đầu vào ổn định. Đặc trưng này không chỉ hỗ trợ việc phân tích xu hướng giá mà còn giúp các mô hình nhận diện các mẫu biến động ngắn hạn, từ đó cải thiện khả năng dự báo trong môi trường thị trường có tính bất ổn cao như tiền mã hóa.
- **Tỷ lệ giá so với trung bình động (Price-MA Ratio):** Được trích xuất để đo lường mức độ lệch của giá hiện tại so với giá trị trung bình của các giá đóng cửa trong một khoảng thời gian gần nhất, thường là 5, 10 hoặc 20 kỳ. Đặc trưng này đóng vai trò quan trọng trong việc xác định các pha thị trường, chẳng hạn như trạng thái quá mua (overbought) khi giá vượt xa trung bình động hoặc quá bán (oversold) khi giá giảm xuống dưới ngưỡng này.
- **Biên độ dao động giá (Price Spread):** Được xác định bằng tỷ lệ giữa độ chênh lệch của giá cao nhất và thấp nhất so với giá đóng cửa trong một khung thời gian nhỏ, chẳng hạn như 5 phút hoặc 15 phút. Đặc trưng này phản ánh mức độ biến động giá trong ngắn hạn, cung cấp thông tin về sự bất ổn của thị trường tại một thời điểm cụ thể. Trong thị trường tiền mã hóa, nơi các biến động giá mạnh thường xảy ra, Price Spread giúp mô hình nhận diện các giai đoạn thị trường có rủi ro cao hoặc các cơ hội giao dịch tiềm năng. Việc tích hợp đặc trưng này vào quá trình dự báo cho phép mô hình học sâu phân tích tốt hơn các mẫu giá ngắn hạn, từ đó cải thiện độ chính xác trong việc dự đoán các xu hướng tức thời.
- **Z-score của khối lượng giao dịch (Volume Z-score):** Z-score của khối lượng giao dịch (Volume Z-score) được tính để đánh giá mức độ bất thường của khối lượng giao dịch hiện tại so với trung bình lịch sử trong một khoảng thời gian gần, thường là 20 hoặc 50 kỳ. Đặc trưng này giúp phát hiện các phiên giao dịch có khối lượng cao bất thường, thường liên quan đến các sự kiện thị trường quan trọng như bùng nổ thanh khoản hoặc trầm lắng bất ngờ.
- **Tỷ lệ khối lượng giao dịch so với trung bình động (Volume-MA Ratio):** Được trích xuất để so sánh sức mua hoặc bán hiện tại với trạng thái bình thường trong khoảng thời gian gần, thường được tính dựa trên trung bình động của khối lượng giao dịch trong 10 hoặc 20 kỳ.
- **Tỷ lệ khối lượng giao dịch so với trung bình động (Volume-MA Ratio)** được trích xuất để so sánh sức mua hoặc bán hiện tại với trạng thái bình thường trong khoảng thời gian gần, thường được tính dựa trên trung bình động của khối lượng giao dịch trong 10 hoặc 20 kỳ.

- **Tính thanh khoản (Liquidity)** – được tính dựa trên tổng giá trị giao dịch quy đổi logarit ( $\text{giá} \times \text{khối lượng}$ ), thể hiện cường độ dòng tiền đổ vào thị trường.
- **Các chỉ báo kỹ thuật** như:
  1. **RSI (Relative Strength Index)**: đo lường sức mạnh xu hướng giá, cho thấy mức độ quá mua hoặc quá bán.
  2. **MACD (Moving Average Convergence Divergence)**: cho biết sự phân kỳ/hội tụ giữa các đường trung bình động, từ đó nhận biết tín hiệu đảo chiều.
  3. **ATR (Average True Range)**: phản ánh độ dao động trung bình của giá, liên quan đến độ rủi ro.
  4. **OBV (On-Balance Volume)**: kết hợp giá và khối lượng để đo dòng tiền tích lũy.
- **Đặc trưng về biến động (Volatility)** – tính bằng độ lệch chuẩn của log return trong các khoảng thời gian ngắn như 30 phút, 1 giờ và 2 giờ. Biến động cao thường đồng nghĩa với rủi ro lớn hơn và khả năng xảy ra biến động mạnh trong giá.
- **Động lượng (Momentum)** – thể hiện xu hướng tăng hoặc giảm của giá trong ngắn và trung hạn bằng cách so sánh trung bình động ngắn hạn và dài hạn. Đây là chỉ số quan trọng trong việc phát hiện tín hiệu đảo chiều.
- **Đặc trưng thời gian** – được trích xuất từ dấu thời gian gồm: giờ trong ngày, thứ trong tuần, xác định phiên cuối tuần và nhận diện phiên giao dịch trong khung giờ thị trường hoạt động mạnh (từ 8h đến 20h). Ngoài ra, các biến này còn được mã hóa tuần hoàn bằng các hàm lượng giác sin và cos, nhằm hỗ trợ mô hình nhận biết tính chu kỳ trong hành vi giao dịch.

Tất cả các đặc trưng nêu trên đều đóng vai trò quan trọng trong việc cung cấp thông tin bổ sung về trạng thái của thị trường, qua đó giúp mô hình nhận diện chính xác hơn các xu hướng và tín hiệu dự báo giá đóng cửa trong tương lai gần.

## 4.2. Chuẩn hóa dữ liệu

Trong tập dữ liệu tiền điện tử, các đặc trưng đầu vào có đơn vị đo lường và độ lớn rất khác nhau. Nếu không chuẩn hóa, mô hình học sâu sẽ bị lệch về những đặc trưng có giá trị tuyệt đối lớn, làm mất cân bằng trong quá trình học. Vì vậy, chuẩn hóa là một bước thiết yếu nhằm đưa các đặc trưng về cùng một thang đo, đảm bảo rằng mỗi đặc trưng có đóng góp công bằng vào quá trình huấn luyện. Trong hệ thống này, hai kỹ thuật chuẩn hóa chính được sử dụng: Min-Max Scaling và Robust Scaling.



### 4.2.1. Min-Max Scaling

#### ♦ Khái niệm

Min-Max Scaling là phương pháp chuẩn hóa tuyến tính, đưa toàn bộ giá trị của một đặc trưng về một khoảng xác định, thường là  $[0,1]$  hoặc  $[-1,1]$ . Phương pháp này giúp dữ liệu có tỷ lệ đóng góp đồng đều và đặc biệt phù hợp với các mô hình học sâu sử dụng hàm kích hoạt có miền xác định như ReLU hoặc tanh.

#### ♦ Công thức

- Giá trị đã chuẩn hóa được tính theo:

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}} \cdot (b - a) + a$$

Trong đó:

$x$ : giá trị gốc của đặc trưng.

$x_{min}$ ,  $x_{max}$ : giá trị nhỏ nhất và lớn nhất của đặc trưng đó trong tập dữ liệu.

$a, b$ : khoảng đích cần chuẩn hóa về (thường là  $a = -1$ ,  $b = 1$ ).

- Khi chuẩn hóa về khoảng  $[0,1]$ , công thức đơn giản hóa còn:

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

#### ♦ Ứng dụng trong hệ thống

- Min-Max Scaling được áp dụng cho các đặc trưng:
- Chỉ báo kỹ thuật: RSI, MACD, ATR, OBV, log return, v.v.
- Đặc trưng thời gian: sin/cos của giờ, ngày trong tuần, cờ nhị phân.

Lý do là các đặc trưng này có phạm vi giá trị ổn định, không có nhiều ngoại lệ lớn nên phù hợp với dạng chuẩn hóa tuyến tính.

### 4.2.2. Robust Scaling

#### ♦ Khái niệm

Robust Scaling là phương pháp chuẩn hóa dựa trên trung vị và khoảng tứ phân vị (IQR) thay vì giá trị lớn nhất/nhỏ nhất. Phương pháp này đặc biệt hiệu quả với dữ liệu chứa nhiều outliers – là những giá trị cực lớn hoặc cực nhỏ bất thường.

Trong tài chính, đặc trưng như giá hoặc khối lượng thường có những đột biến lớn do sự kiện thị trường hoặc lỗi ghi nhận, vì vậy không thể dùng Min-Max Scaling, vì các giá trị cực trị sẽ kéo dẫn toàn bộ thang đo.

#### ♦ Công thức

- Giá trị đã chuẩn hóa được tính như sau:

$$x_{scaled} = \frac{x - median(x)}{IQR(x)}$$

Trong đó:

median(x): trung vị của đặc trưng.

IQR(x) = Q3–Q1: là khoảng tứ phân vị, với:

Q1: phân vị thứ 25%

Q3: phân vị thứ 75%

Ứng dụng trong hệ thống:

- Robust Scaling được áp dụng cho:
  - Đặc trưng về giá: open, high, low, close, price\_ma\_ratio, price\_spread
  - Đặc trưng về khối lượng: volume, volume\_ma\_ratio, volume\_zscore, liquidity
- Lý do là các đặc trưng này dễ có đột biến lớn (ví dụ: pump & dump, khối lượng giao dịch bất thường), và trung vị kết hợp với IQR sẽ giúp chuẩn hóa ổn định hơn nhiều so với việc dựa vào giá trị lớn nhất/nhỏ nhất.

#### 4.2.3. Phân nhóm đặc trưng theo phương pháp chuẩn hóa

Nhóm đặc trưng	Các cột liên quan	Kỹ thuật chuẩn hóa
<b>Giá cả (Price)</b>	open, high, low, close, price_ma_ratio, price_spread	Robust Scaling
<b>Khối lượng (Volume)</b>	volume, volume_ma_ratio, volume_zscore, liquidity	Robust Scaling
<b>Chỉ báo kỹ thuật</b>	RSI, MACD, ATR, OBV, log_returns, volatility, momentum	Min-Max Scaling [-1, 1]
<b>Đặc trưng thời gian</b>	sin(hour), cos(hour), sin(dow), cos(dow), is_weekend, market_open	Min-Max Scaling [0, 1]

#### Lợi ích của việc phân chia chuẩn hóa theo nhóm

- Tăng tính ổn định của mô hình: mỗi nhóm dữ liệu được xử lý phù hợp với đặc điểm phân phối của nó.
- Tránh hiện tượng exploding gradient trong huấn luyện mạng nơ-ron khi đặc trưng chưa được chuẩn hóa đúng cách.
- Tối ưu hóa tốc độ hội tụ của các thuật toán tối ưu như Adam hoặc SGD.

#### 4.3. Tăng cường dữ liệu (Data Augmentation)

Trong các bài toán học sâu với dữ liệu chuỗi thời gian tài chính, một trong những thách thức lớn nhất là hiện tượng **quá khớp (overfitting)**, tức là mô hình học thuộc quá sát vào dữ liệu huấn luyện và không còn khả năng khái quát hóa tốt trên dữ liệu chưa thấy. Điều này thường xảy ra khi mô hình quá phức tạp trong khi dữ liệu đầu vào không đủ đa dạng, hoặc có các quy luật cứng nhắc bị học sai.

Để khắc phục tình trạng này, hệ thống áp dụng một tập hợp các kỹ thuật **tăng cường dữ liệu (data augmentation)** nhằm mô phỏng nhiều kịch bản thị trường khác nhau. Việc tạo ra các biến thể nhân tạo của dữ liệu không chỉ làm tăng kích thước tập huấn luyện một cách gián tiếp mà còn giúp mô hình trở nên "cứng cáp", không bị phụ thuộc vào một số mẫu cụ thể hoặc nhiều thị trường tạm thời.

Tăng cường dữ liệu được áp dụng ngẫu nhiên tại mỗi bước lấy mẫu trong quá trình huấn luyện, với các kỹ thuật sau:

#### 4.3.1. Che cục bộ (Local Mean Masking)

Phương pháp này mô phỏng hiện tượng thiếu dữ liệu cục bộ, thường gặp trong dữ liệu thực tế do lỗi ghi nhận, mất kết nối hoặc các yếu tố kỹ thuật.

Cụ thể, một đoạn ngắn liên tiếp trong chuỗi thời gian (ví dụ 5 đến 10 bước liên tục) sẽ bị thay thế bằng **giá trị trung bình của vùng lân cận mở rộng**, tức là không chỉ lấy trung bình trong đoạn bị che mà còn lấy thêm dữ liệu hai bên.

Phương pháp này buộc mô hình phải suy luận thông tin từ bối cảnh xung quanh thay vì chỉ dựa vào giá trị tuyệt đối, từ đó nâng cao khả năng khái quát hóa khi gặp các đoạn dữ liệu thiếu trong thực tế.

#### 4.3.2. Thêm nhiễu Gaussian (Adaptive Gaussian Noise)

Một kỹ thuật cơ bản nhưng hiệu quả để làm cho mô hình ít nhạy cảm với biến động nhỏ trong dữ liệu là **thêm nhiễu ngẫu nhiên phân phối chuẩn (Gaussian noise)** vào chuỗi đầu vào.

Nhiễu được sinh ra từ phân phối chuẩn có trung bình bằng 0 và độ lệch chuẩn tỷ lệ với độ biến động của đặc trưng:

$$x'_t = x_t + \varepsilon_t \quad \text{với} \quad \varepsilon_t \sim N(0, \sigma^2)$$

Trong đó  $\sigma$  được điều chỉnh theo **tiến độ huấn luyện**, thường tăng dần theo số epoch nhằm giúp mô hình dần thích nghi với thị trường ngày càng nhiễu hơn. Kỹ thuật này đặc biệt hữu ích trong các giai đoạn đầu huấn luyện, giúp tránh tối ưu hóa cứng nhắc theo một đường giá cụ thể.

#### 4.3.3. Tỷ lệ ngẫu nhiên (Smart Scaling)

Phương pháp này nhằm mô phỏng sự thay đổi quy mô dữ liệu trong thị trường thực tế. Thay vì thay đổi toàn bộ chuỗi, hệ thống chọn ngẫu nhiên một số đặc trưng đầu vào (chẳng hạn như

volume hoặc liquidity) và nhân chúng với một hệ số biến thiên nhỏ, lấy từ phân phối đồng đều trong khoảng  $[0.9, 1.1]$ .

Mặc dù sự thay đổi này là nhỏ, nhưng nó giúp mô hình **học được tính bất định trong hành vi thị trường**, nơi mà các chỉ số có thể dao động nhẹ do nhiễu, chênh lệch sàn giao dịch, hoặc thay đổi theo điều kiện thị trường.

Smart Scaling giữ nguyên mối quan hệ tương đối giữa các bước thời gian trong chuỗi, tức là vẫn giữ được "hình dạng" dữ liệu gốc, nhưng giúp mô hình không học lệ thuộc vào tuyệt đối các giá trị cụ thể.

#### 4.3.4. Giãn thời gian (Time Warping)

Trong thực tế, tốc độ thay đổi của thị trường không phải lúc nào cũng đều đặn: có lúc biến động nhanh, có lúc chậm. Để mô phỏng điều này, chuỗi thời gian đầu vào được **co hoặc giãn tuyến tính theo thời gian** thông qua phép nội suy (interpolation).

Chuỗi gốc sẽ được kéo dài hoặc rút ngắn bằng một hệ số ngẫu nhiên trong khoảng (ví dụ: từ 0.8 đến 1.2), sau đó được cắt hoặc đệm lại đúng độ dài ban đầu. Điều này buộc mô hình phải học các đặc trưng không phụ thuộc vào tốc độ biến động cố định, giúp tăng khả năng thích nghi với các điều kiện thị trường linh hoạt hơn.

Time Warping đặc biệt hiệu quả trong việc mô phỏng các tình huống "shock" thị trường hoặc "sideways" kéo dài.

#### 4.3.5. Bỏ đặc trưng ngẫu nhiên (Feature Dropout)

Mô hình học sâu thường có xu hướng phụ thuộc vào một số đặc trưng nhất định. Nếu các đặc trưng đó không khả dụng trong tương lai (do nhiễu, mất dữ liệu...), mô hình sẽ mất khả năng dự đoán.

Để tránh điều đó, trong mỗi lần lấy mẫu, hệ thống sẽ **ngẫu nhiên vô hiệu hóa (gán bằng 0)** một số đặc trưng đầu vào với xác suất nhất định. Quá trình này tương tự như kỹ thuật Dropout trong huấn luyện mạng nơ-ron, nhưng áp dụng ở cấp độ đặc trưng thay vì node.

Nhờ vậy, mô hình học cách **tận dụng nhiều nguồn tín hiệu khác nhau**, không quá phụ thuộc vào một số đặc trưng nổi bật, từ đó tăng tính bền vững (robustness) khi gặp dữ liệu thực tế không hoàn hảo.

## 4.4. Các Mô hình Dự đoán:

Để giải quyết bài toán dự đoán chuỗi thời gian trong lĩnh vực tài chính, đặc biệt là dữ liệu giá tiền mã hóa có tính biến động mạnh, nhóm nghiên cứu đã xây dựng và triển khai ba kiến trúc mô hình học sâu: **mô hình LSTM thuần**, **mô hình CNN-LSTM**, và **mô hình CNN-LSTM kết hợp Attention**. Mỗi mô hình thể hiện một hướng tiếp cận khác nhau trong việc trích xuất đặc trưng và mô hình hóa phụ thuộc theo thời gian. Các mô hình được thiết kế theo cấu trúc thống nhất, nhưng mỗi mô hình sẽ được mở rộng theo chiều sâu hoặc tích hợp thêm các khối xử lý đặc biệt để tăng cường khả năng học của mạng nơ-ron.

### 4.4.1. Mạng LSTM thuần túy

Long Short-Term Memory (LSTM) là một dạng mạng hồi tiếp cải tiến, được thiết kế để ghi nhớ các quan hệ dài hạn trong chuỗi dữ liệu thông qua cơ chế bộ nhớ và điều khiển thông tin bằng ba cổng: cổng quên, cổng đầu vào và cổng đầu ra. Mỗi đơn vị LSTM cho phép mô hình duy trì thông tin qua nhiều bước thời gian mà không gặp hiện tượng gradient biến mất.

Trong cấu trúc được đề xuất, đầu vào bao gồm một chuỗi thời gian nhiều chiều. Chuỗi này được đưa vào khối xử lý sơ cấp, nơi dữ liệu được ánh xạ sang không gian ẩn thông qua các phép biến đổi tuyến tính và phi tuyến, kèm theo chuẩn hóa lớp và dropout nhằm ổn định huấn luyện và tăng cường khả năng khái quát. Sau đó, chuỗi được đưa qua một tầng LSTM duy nhất, hoạt động theo hướng tiến (unidirectional), tạo ra một dãy các vector biểu diễn ẩn cho từng bước thời gian. Thay vì sử dụng toàn bộ dãy đầu ra, mô hình chỉ lấy biểu diễn tại bước thời gian cuối cùng làm đặc trưng đại diện cho toàn bộ chuỗi đầu vào.

Lựa chọn này xuất phát từ giả định rằng trạng thái cuối cùng của bộ nhớ LSTM đã tích lũy đủ thông tin phản ánh toàn bộ chuỗi lịch sử, từ đó có thể sử dụng làm cơ sở cho việc dự đoán. Biểu diễn cuối cùng được đưa vào khối mạng tuyến tính để sinh ra chuỗi đầu ra tương ứng với khoảng thời gian dự đoán tiếp theo.

Cấu trúc này có tính đơn giản, phù hợp cho các bài toán chuỗi có quy luật dài hạn rõ ràng hoặc dữ liệu đã được xử lý tốt để giảm nhiễu.

### 4.4.2. Mô hình CNN-LSTM

Đối với các chuỗi dữ liệu tài chính thực tế, đặc trưng cục bộ và dao động ngắn hạn thường đóng vai trò quan trọng nhưng khó phát hiện nếu chỉ dựa vào cơ chế tuần tự của LSTM. Để giải quyết vấn đề này, kiến trúc kết hợp giữa mạng tích chập một chiều (1D Convolutional Neural Network – CNN) và LSTM đã được xây dựng.

Thành phần CNN được đặt trước tầng LSTM với mục đích trích xuất đặc trưng ngắn hạn và các mẫu dao động có tính lặp lại trong chuỗi. CNN sử dụng các bộ lọc với kích thước kernel nhỏ để quét trên toàn chuỗi đầu vào, nhằm phát hiện các đặc trưng cục bộ mà các lớp tuần tự không dễ nhận diện. Tầng tích chập này cũng góp phần làm giảm nhiễu và làm nổi bật các yếu tố quan trọng, giúp đơn giản hóa đầu vào cho LSTM.

Sau khi đi qua CNN, chuỗi được chuyển đổi về dạng chuẩn để đưa vào tầng LSTM, nơi học các quan hệ chuỗi dài hạn. LSTM trong kiến trúc này giữ nguyên nguyên lý hoạt động như mô hình trước, tuy nhiên đầu vào của nó đã được cải tiến nhờ quá trình rút trích đặc trưng trước đó. Tương tự, biểu diễn tại bước thời gian cuối của LSTM được chọn làm đặc trưng tổng hợp và đưa vào mạng con đầu ra để sinh dự đoán.

Sự kết hợp này cho phép mô hình đồng thời tận dụng khả năng phát hiện tín hiệu ngắn hạn nhờ CNN và khả năng ghi nhớ dài hạn từ LSTM. Cấu trúc này đặc biệt hữu ích trong các tình huống mà chuỗi thời gian chứa nhiều biến động cục bộ và các sự kiện có tính ngữ cảnh ngắn hạn.

#### 4.4.3. Mô hình CNN-LSTM Tích hợp Attention

Mặc dù mô hình CNN-LSTM đã thể hiện khả năng học đa cấp thông tin chuỗi, một hạn chế còn tồn tại là việc lựa chọn biểu diễn tại bước thời gian cuối cùng làm đặc trưng đại diện có thể bỏ qua nhiều tín hiệu quan trọng xuất hiện trước đó trong chuỗi. Để khắc phục điểm yếu này, nhóm nghiên cứu đề xuất bổ sung **cơ chế chú ý (Attention mechanism)** nhằm giúp mô hình tự động học trọng số cho từng bước thời gian và tập trung vào những phần thông tin có giá trị nhất cho nhiệm vụ dự đoán.

Trong kiến trúc này, dữ liệu được xử lý tuần tự qua khối CNN và LSTM như trong mô hình trước. Tuy nhiên, thay vì chỉ sử dụng bước thời gian cuối cùng từ LSTM, toàn bộ dãy đầu ra từ LSTM được giữ lại và đưa vào khối Attention. Khối này bao gồm một chuỗi các phép biến đổi tuyến tính và phi tuyến nhằm ánh xạ từng bước thời gian về một không gian trọng số chú ý. Hàm softmax được sử dụng để chuẩn hóa trọng số thành phân phối xác suất, từ đó tính toán mức độ đóng góp tương đối của mỗi bước thời gian.

Trọng số chú ý sau đó được dùng để thực hiện phép tổng có trọng số lên chuỗi đầu ra của LSTM, tạo thành một vector ngữ cảnh duy nhất. Vector này có tính chất tổng hợp nhưng nhấn mạnh vào các thời điểm có ảnh hưởng lớn đến kết quả dự đoán. Cuối cùng, vector ngữ cảnh được đưa vào tầng tuyến tính đầu ra để tạo thành chuỗi dự đoán tương lai.

Cơ chế Attention đóng vai trò như một bộ lọc thông minh, cho phép mô hình không chỉ “nhìn” vào toàn bộ chuỗi mà còn biết “tập trung” vào phần nào. Điều này mang lại hiệu quả cao trong bối

cảnh dữ liệu phi tuyến, bất ổn định và chứa nhiều nhiễu như giá tiền mã hóa, nơi mà một vài bước thời gian cụ thể có thể chỉ phối toàn bộ xu hướng tiếp theo.

Việc triển khai ba mô hình với cấu trúc tăng dần về khả năng biểu diễn cho phép đánh giá rõ ràng vai trò và đóng góp của từng thành phần trong hiệu quả dự đoán. Mô hình cuối cùng (CNN-LSTM-Attention) thể hiện tiềm năng cao nhất trong việc xử lý dữ liệu chuỗi phi tuyến, bất ổn định và có tính ngữ cảnh phức tạp như dữ liệu tài chính.

## Chương 5. KẾT QUẢ THỰC NGHIỆM VÀ THẢO LUẬN

Kết quả Thực nghiệm và Thảo luận có chức năng chính là trình bày và phân tích kết quả từ quá trình thực hiện nghiên cứu. Chương đóng vai trò quan trọng trong việc đánh giá hiệu suất của mô hình LSTM và so sánh kết quả với các phương pháp truyền thống. Ngoài ra, chương cũng tập trung vào việc giải thích ý nghĩa của các kết quả và mối quan hệ giữa các biến quan trọng.

### 5.1. Các Mô hình Dự đoán:

Trong chương này, chúng tôi trình bày quy trình xây dựng mô hình, tổ chức dữ liệu, môi trường huấn luyện và các kỹ thuật được áp dụng trong quá trình thực nghiệm. Nghiên cứu tập trung vào ba kiến trúc học sâu: LSTM thuần túy, CNN-LSTM, và CNN-LSTM tích hợp Attention. Mỗi mô hình được thiết kế nhằm khai thác đặc trưng chuỗi thời gian tài chính – vốn thường chứa nhiều biến động ngắn hạn, xu hướng không ổn định, và các đột biến khó dự đoán.

#### 5.1.1. Dữ liệu và cấu trúc chuỗi đầu vào

Dữ liệu đầu vào được thu thập từ thị trường tiền mã hóa với độ phân giải 5 phút, sử dụng giá **đóng cửa (close)** làm biến mục tiêu. Bên cạnh đó, tập dữ liệu cũng đã được tiền xử lý, chuẩn hóa và phân nhóm đặc trưng theo phương pháp phù hợp (Min-Max hoặc Robust Scaling) như đã trình bày trong các phần trước.

Trong giai đoạn thực nghiệm, dữ liệu được tổ chức thành các cặp chuỗi đầu vào – đầu ra theo kỹ thuật sliding window:

- **Chuỗi đầu vào:** bao gồm 288 bước thời gian gần nhất, tương ứng với 24 giờ giao dịch
- **Chuỗi đầu ra:** là 36 bước tiếp theo, tương ứng với khoảng 3 giờ dự báo

Mỗi mẫu huấn luyện phản ánh một khung thời gian giao dịch liên tục, cho phép mô hình học được cả ngữ cảnh cận thời gian lẫn xu hướng dài hơn.



### 5.1.2. Mô hình huấn luyện

Trong quá trình thực nghiệm, ba mô hình học sâu khác nhau đã được triển khai để đánh giá khả năng dự đoán chuỗi thời gian giá đóng cửa tiền mã hóa. Các mô hình đều sử dụng cùng một đầu vào có cấu trúc chuỗi gồm 288 bước thời gian và dự đoán chuỗi đầu ra gồm 36 bước tiếp theo. Tất cả các mô hình đều được xây dựng bằng thư viện PyTorch, với kiến trúc được thiết kế và điều chỉnh nhằm tối ưu hiệu quả học biểu diễn theo thời gian. Dưới đây là mô tả chi tiết từng mô hình:

#### a) Mô hình LSTM thuần túy

Đây là mô hình học chuỗi thời gian cơ bản, sử dụng một mạng LSTM để học trực tiếp từ chuỗi giá đóng cửa và đặc trưng kỹ thuật. Mô hình có cấu trúc như sau:

- Số đặc trưng đầu vào (enc\_in): 26 đặc trưng tại mỗi bước thời gian
- Độ dài chuỗi đầu vào (seq\_len): 288 bước (tương đương 24 giờ)
- Độ dài chuỗi đầu ra (pred\_len): 36 bước (tương đương 3 giờ)
- Chiều ẩn LSTM (d\_model): 128
- Số lớp LSTM (e\_layers): 1
- Số "attention heads": 4 (chưa sử dụng trong mô hình này nhưng có sẵn cho mô hình mở rộng)
- Dropout: 0.5
- Kích thước đầu ra: 1 giá trị liên tục cho mỗi bước dự đoán

Mô hình sử dụng một tầng LSTM một chiều, không có cơ chế chú ý hay tích chập, và biểu diễn tại bước thời gian cuối cùng được trích xuất làm đặc trưng tổng hợp để đưa vào tầng tuyến tính đầu ra.

#### b) Mô hình CNN-LSTM

Để cải thiện khả năng học đặc trưng ngắn hạn, mô hình CNN-LSTM bổ sung thêm một tầng mạng tích chập 1 chiều (1D Convolutional) trước tầng LSTM. Cấu trúc cụ thể:

- Số đặc trưng đầu vào (enc\_in): 26
- Chiều dài chuỗi đầu vào (seq\_len): 288

- Chiều dài chuỗi đầu ra (pred\_len): 36
- Chiều ẩn của LSTM (d\_model): 128
- Số filters (cnn\_out\_channels): 64
- Dropout: 0.4
- Kích thước đầu ra: 1 giá trị tại mỗi bước

Tầng CNN sử dụng kernel có kích thước nhỏ để trích xuất các mẫu ngắn hạn (short-term patterns). Sau CNN, dữ liệu được định hình lại và đưa vào LSTM, nơi mô hình học các phụ thuộc theo thời gian. Biểu diễn tại bước thời gian cuối cùng tiếp tục được đưa qua mạng đầu ra để dự đoán chuỗi 36 bước tiếp theo.

### c) Mô hình CNN-LSTM Attention

Đây là mô hình phức tạp nhất trong ba kiến trúc, được mở rộng từ mô hình CNN-LSTM bằng cách tích hợp thêm cơ chế Attention sau tầng LSTM. Thay vì chỉ sử dụng biểu diễn tại bước thời gian cuối cùng, mô hình sử dụng toàn bộ chuỗi đầu ra từ LSTM và học phân phối trọng số cho từng bước thời gian, từ đó tạo ra vector ngữ cảnh tổng hợp có trọng số.

Cấu trúc chi tiết:

- Số đặc trưng đầu vào (enc\_in): 26
- Chiều dài chuỗi đầu vào (seq\_len): 288
- Chiều dài chuỗi đầu ra (pred\_len): 36
- Chiều ẩn của LSTM (d\_model): 128
- Số filters trong CNN (cnn\_out\_channels): 64
- Dropout: 0.4
- Kích thước đầu ra: 1 giá trị mỗi bước

Khối Attention sử dụng một tầng tuyến tính để tính điểm mức độ quan trọng cho từng bước thời gian trong chuỗi đầu ra của LSTM, sau đó sử dụng softmax để tạo ra trọng số phân bố. Toàn bộ chuỗi được tổng hợp có trọng số và đưa vào mạng dự đoán.

Mô hình này có khả năng chọn lọc và tập trung vào các thời điểm có ảnh hưởng lớn đến xu hướng giá, đồng thời vẫn giữ được năng lực học các mẫu ngắn và dài hạn thông qua CNN và LSTM.

### 5.1.3. Quy trình huấn luyện

Tất cả các mô hình được huấn luyện theo một quy trình thống nhất với mục tiêu đảm bảo tính công bằng trong so sánh hiệu quả giữa các kiến trúc khác nhau. Quá trình huấn luyện được tối ưu hóa về cả hiệu năng tính toán lẫn hiệu suất hội tụ, đồng thời tích hợp nhiều kỹ thuật hiện đại trong đào tạo mô hình học sâu.

#### Chiến lược tối ưu hóa:

Các mô hình sử dụng **thuật toán tối ưu hóa AdamW** với hệ số học khởi tạo thấp nhằm đảm bảo ổn định trong giai đoạn đầu huấn luyện. Để điều chỉnh tốc độ học linh hoạt theo tiến trình huấn luyện, chúng tôi sử dụng **lịch trình giảm tốc độ học (learning rate scheduler) theo dạng Cosine với warmup**.

Trong giai đoạn warmup, tốc độ học được tăng dần từ 0 đến giá trị tối đa đã định sẵn. Sau đó, theo tiến trình huấn luyện, learning rate được giảm dần theo dạng hình sin về gần 0. Điều này giúp mô hình có khởi đầu ổn định, đồng thời tránh rơi vào local minima quá sớm trong các epoch đầu tiên.

#### Thiết lập huấn luyện

- **Hàm mất mát:** MSELoss (Mean Squared Error), phù hợp với bài toán hồi quy giá trị liên tục
- **Thuật toán tối ưu:** AdamW với weight decay nhẹ ( $1e-3$ )
- **Batch size:** 128
- **Epoch tối đa:** 100
- **Learning rate khởi tạo:** giá trị được cấu hình qua file YAML
- **Warmup epochs:** 5 epoch đầu

#### Tăng tốc và ổn định huấn luyện

- **Mixed Precision (AMP):** Áp dụng tự động với torch.cuda.amp, giúp giảm sử dụng bộ nhớ và tăng tốc độ huấn luyện mà không ảnh hưởng đáng kể đến độ chính xác

- **Gradient Clipping:** Được sử dụng để ổn định huấn luyện, giới hạn giá trị gradient ở mức tối đa nhằm tránh hiện tượng nổ gradient

#### **Chiến lược dừng sớm (*Early Stopping*):**

Một cơ chế kiểm soát overfitting được áp dụng thông qua kỹ thuật **Early Stopping**, theo dõi giá trị hàm mất mát trên tập huấn luyện. Nếu mô hình không còn cải thiện sau 30 vòng lặp liên tiếp, quá trình huấn luyện sẽ tự động dừng lại. Việc này giúp tiết kiệm tài nguyên và tránh mô hình học quá mức vào dữ liệu huấn luyện.

#### **5.1.4. Phân chia dữ liệu**

Dữ liệu được phân chia theo thời gian thực tế nhằm mô phỏng đúng điều kiện của một hệ thống dự đoán trong môi trường vận hành:

- **Tập huấn luyện:** chiếm 80% đầu chuỗi thời gian
- **Tập val:** là 20% còn lại, sử dụng để đánh giá mô hình sau huấn luyện 1 epoch

Không thực hiện xáo trộn (shuffle) dữ liệu, nhằm duy trì cấu trúc thời gian và đảm bảo không có hiện tượng rò rỉ thông tin từ tương lai về quá khứ.

#### **5.1.5. Môi trường thực nghiệm**

Toàn bộ quá trình huấn luyện được thực hiện trên nền tảng **Kaggle Notebooks**, với cấu hình phần cứng và phần mềm như sau:

- **GPU:** NVIDIA Tesla P100 (16 GB bộ nhớ GPU)
- **RAM hệ thống:** 13 GB
- **Thư viện lập trình:**
  - PyTorch (huấn luyện mô hình)
  - pandas, numpy (xử lý dữ liệu)
  - scikit-learn (chuẩn hóa, chia dữ liệu, đánh giá)
  - matplotlib, seaborn (trực quan hóa kết quả)

## 5.2. Tiêu chí đánh giá

Để đánh giá hiệu quả dự báo của các mô hình được xây dựng trong nghiên cứu, chúng tôi sử dụng một tập hợp các độ đo định lượng phổ biến trong lĩnh vực học máy và phân tích chuỗi thời gian. Cụ thể, các độ đo bao gồm: **MAE** (Mean Absolute Error), **MAPE** (Mean Absolute Percentage Error), **SMAPE** (Symmetric Mean Absolute Percentage Error), **MSE** (Mean Squared Error), **RMSE** (Root Mean Square Error), **R<sup>2</sup>** (Hệ số xác định). Chúng đo lường mức độ sai lệch giữa các giá trị dự đoán và giá trị thực tế. Với  $n$  là số mẫu trong tập dữ liệu,  $Y_t$  biểu thị giá trị thực tế của mẫu thứ  $t$ , và  $\hat{Y}_t$  biểu thị giá trị dự đoán của mẫu thứ  $t$ .

### MAPE (Mean Absolute Percentage Error)

MAPE là một chỉ số đánh giá độ chính xác của mô hình dự báo, thể hiện sai số **trung bình** giữa giá trị dự đoán và giá trị thực tế dưới dạng **phần trăm**.

$$MAPE = \frac{1}{n} \sum_{t=1}^n \left| \frac{Y_t - \hat{Y}_t}{Y_t} \right|$$

- Dễ diễn giải trong thực tế (ví dụ: "mô hình sai khoảng 8% so với giá trị thực").
- Tuy nhiên, MAPE **không ổn định** khi  $Y_t$  gần 0 (chia cho số rất nhỏ dẫn đến giá trị rất lớn hoặc vô cùng).

### MAE (Mean Absolute Error)

MAE đo lường sự sai lệch trung bình giữa các giá trị dự đoán và giá trị thực tế, sử dụng giá trị tuyệt đối trong quá trình tính toán. Giá trị MAE thấp hơn biểu thị độ chính xác mô hình tốt hơn.

$$MAE = \frac{1}{n} \sum_{t=1}^n |Y_t - \hat{Y}_t|$$

- **Giá trị nhỏ** hơn thể hiện mô hình dự đoán gần sát với thực tế hơn.
- **Không nhạy cảm** với các sai số lớn bất thường (outlier) như MSE hay RMSE.

### RMSE (Root Mean Squared Error)

RMSE đánh giá độ chính xác của mô hình dự đoán so với dữ liệu thực. Giá trị RMSE thấp hơn cho thấy mức độ chính xác của mô hình cao hơn.

$$RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^n (Y_t - \hat{Y}_t)^2}$$

- RMSE thường có giá trị lớn hơn MAE.

- Rất hữu ích khi chúng tôi muốn nhấn mạnh mức độ ảnh hưởng của **outlier** đến mô hình.

## R<sup>2</sup> Score (R Squared Score)

R<sup>2</sup>-score sẽ đánh giá mức độ mô hình phù hợp với dữ liệu. Nó được tính bằng cách lấy tổng của bình phương của sự chênh lệch giữa các giá trị dự đoán và giá trị thực tế, chia cho tổng của bình phương của sự chênh lệch giữa giá trị trung bình của giá trị thực tế và giá trị thực tế, và trừ kết quả này từ 1. Điểm R<sup>2</sup> cao hơn biểu thị sự phù hợp tốt hơn của mô hình với dữ liệu, đó là một chỉ số đánh giá độ chính xác của mô hình.

$$R^2 = 1 - \frac{\sum_{t=1}^n (Y_t - \hat{Y}_t)^2}{\sum_{t=1}^n (Y_t - \bar{Y})^2}$$

- R<sup>2</sup> gần 1: mô hình dự đoán tốt.
- R<sup>2</sup> gần 0: mô hình gần như không cải thiện gì so với chỉ đoán trung bình.

## MSE (Mean Squared Error – Sai số bình phương trung bình)

MSE là một độ đo phổ biến, tính trung bình bình phương sai số giữa dự đoán và thực tế.

$$MSE = \frac{1}{n} \sum_{t=1}^n (Y_t - \hat{Y}_t)^2$$

- Nhấn mạnh các sai số lớn do việc **bình phương** sai số.
- Phù hợp khi **sai số lớn cần được trừng phạt mạnh hơn** trong bài toán.

## SMAPE (Symmetric Mean Absolute Percentage Error)

SMAPE là một phiên bản cải tiến của MAPE, giúp khắc phục nhược điểm của MAPE khi  $Y_t$  gần 0 bằng cách chia cho tổng tuyệt đối của giá trị thực và giá trị dự đoán.

$$SMAPE = \frac{100}{n} \sum_{t=1}^n \frac{2 \cdot |Y_t - \hat{Y}_t|}{|Y_t| + |\hat{Y}_t| + \epsilon}$$

- Giúp cân bằng sai số giữa trường hợp **dự đoán thấp** và **dự đoán cao**.
- Rất phù hợp trong các bài toán tài chính, năng lượng, hay dự báo nhu cầu, nơi dữ liệu có thể dao động mạnh.

## 5.3. Phân tích kết quả thực nghiệm

Trong nghiên cứu này, chúng tôi tiến hành huấn luyện mô hình dự đoán giá Bitcoin trên tập dữ liệu lịch sử, sau đó đánh giá hiệu suất mô hình bằng cách kiểm thử trên tập dữ liệu giá Bitcoin 30 ngày gần nhất được lấy từ sàn Binance, với tần suất lấy mẫu 5 phút. Ba kiến trúc được triển khai

bao gồm: mô hình LSTM thuần, mô hình kết hợp CNN-LSTM, và mô hình tích hợp Attention vào LSTM. Việc đánh giá hiệu quả dự đoán được thực hiện thông qua bốn chỉ số chính: **MSE (Mean Squared Error)**, **MAE (Mean Absolute Error)**, **RMSE (Root Mean Squared Error)** và **R<sup>2</sup> (Coefficient of Determination)**.

Trong nghiên cứu này, ba mô hình học sâu bao gồm **LSTM**, **CNN-LSTM**, và **CNN-LSTM kết hợp Attention** đã được xây dựng và huấn luyện nhằm giải quyết bài toán dự báo giá tiền điện tử trong tương lai gần. Tập dữ liệu huấn luyện được lấy từ dữ liệu giá lịch sử của đồng **Bitcoin**, trong khi đó, tập kiểm thử được thu thập độc lập từ sàn giao dịch **Binance**, với **các điểm dữ liệu có tần suất 5 phút**, trải dài trong **30 ngày gần nhất** so với thời điểm thực hiện thực nghiệm.

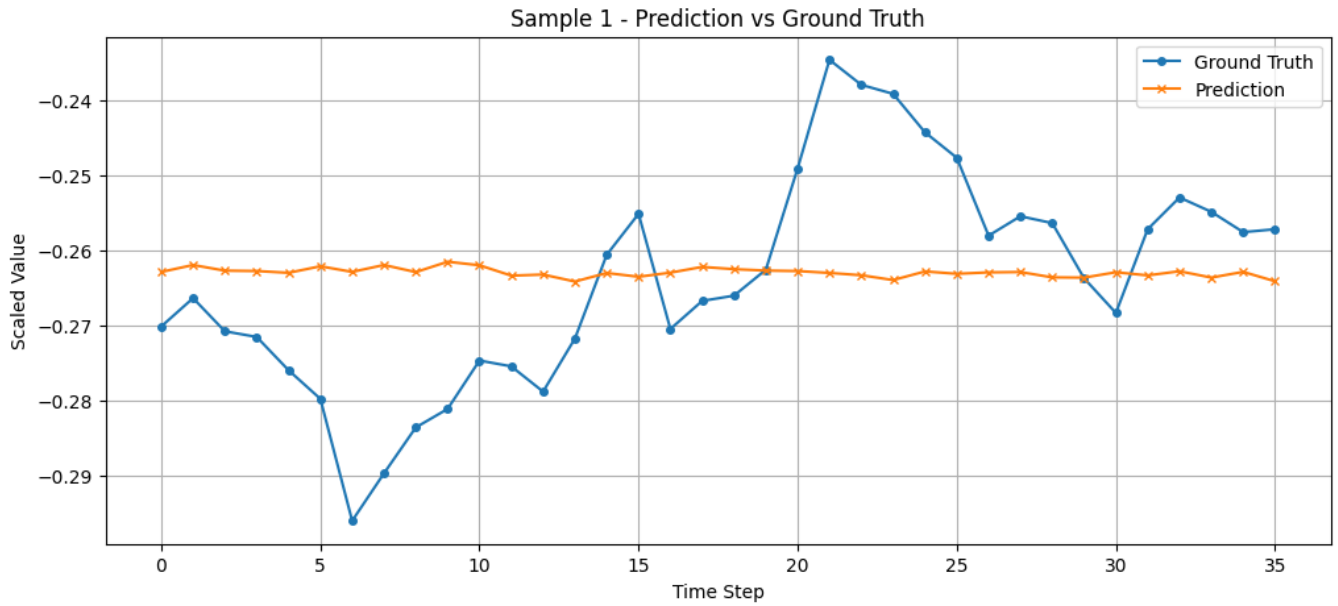
### 5.3.1 Hiệu suất mô hình LSTM

MSE	MAE	RMSE	R <sup>2</sup>
0.0109	0.0755	0.1043	0.9677

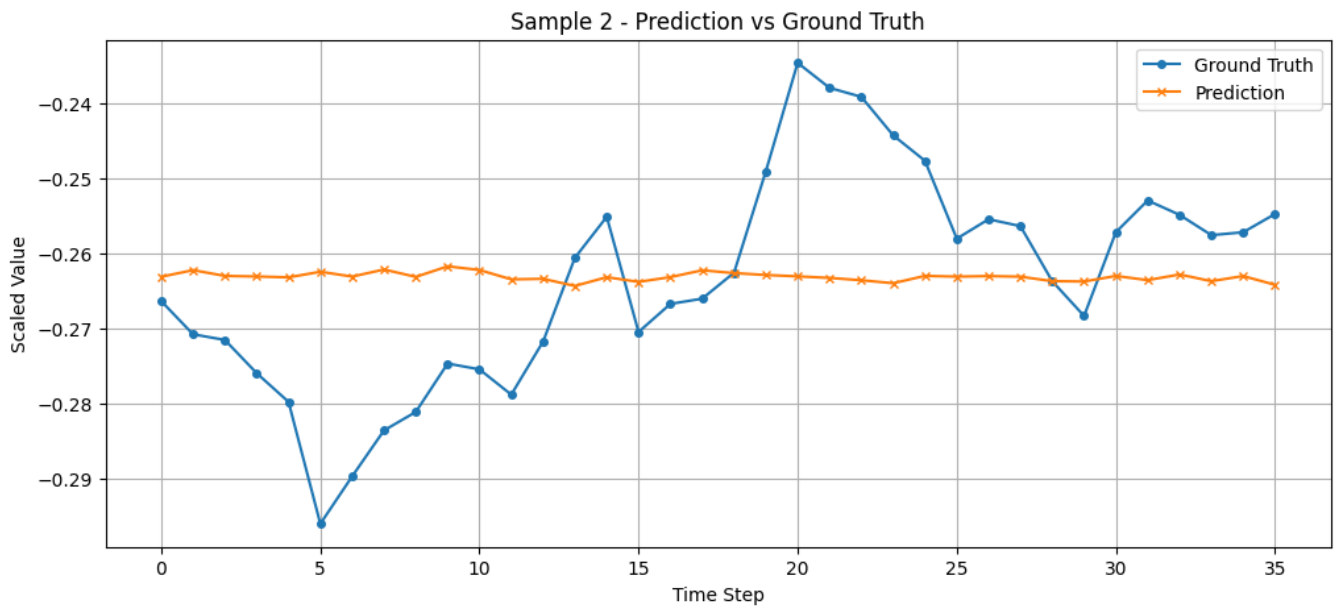
Mô hình Long Short-Term Memory (LSTM), một kiến trúc mạng nơ-ron hồi tiếp được thiết kế chuyên biệt để xử lý và học các phụ thuộc dài hạn trong chuỗi thời gian, đã thể hiện hiệu suất vượt trội trong thực nghiệm. Với sai số tuyệt đối trung bình (MAE) là 0.0755 và sai số phần trăm trung bình (MSE) chỉ 0.0109 mô hình này cho thấy khả năng dự báo chính xác giá trị tương lai trong chuỗi giá có độ biến động cao như tiền điện tử.

Hơn nữa, giá trị RMSE đạt 2.7940 và hệ số xác định  $R^2 = 0.9677$  phản ánh rằng mô hình có thể giải thích tới 96.77% phương sai trong dữ liệu kiểm thử. Đây là một chỉ báo mạnh mẽ về mức độ phù hợp giữa giá trị dự báo và giá trị thực tế, cũng như khả năng nắm bắt được xu hướng biến động chính của chuỗi giá.

**Biểu đồ dự đoán so với giá trị thực (5 mẫu ngẫu nhiên):**

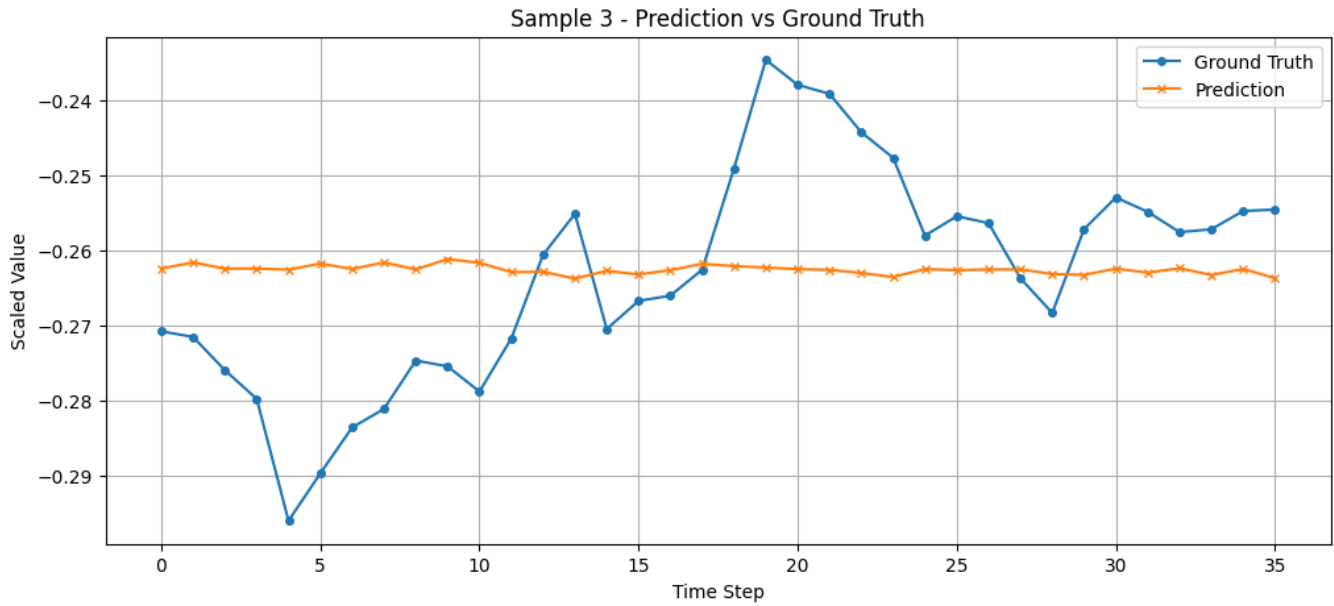


Hình 5.1. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình LSTM – Mẫu 1.

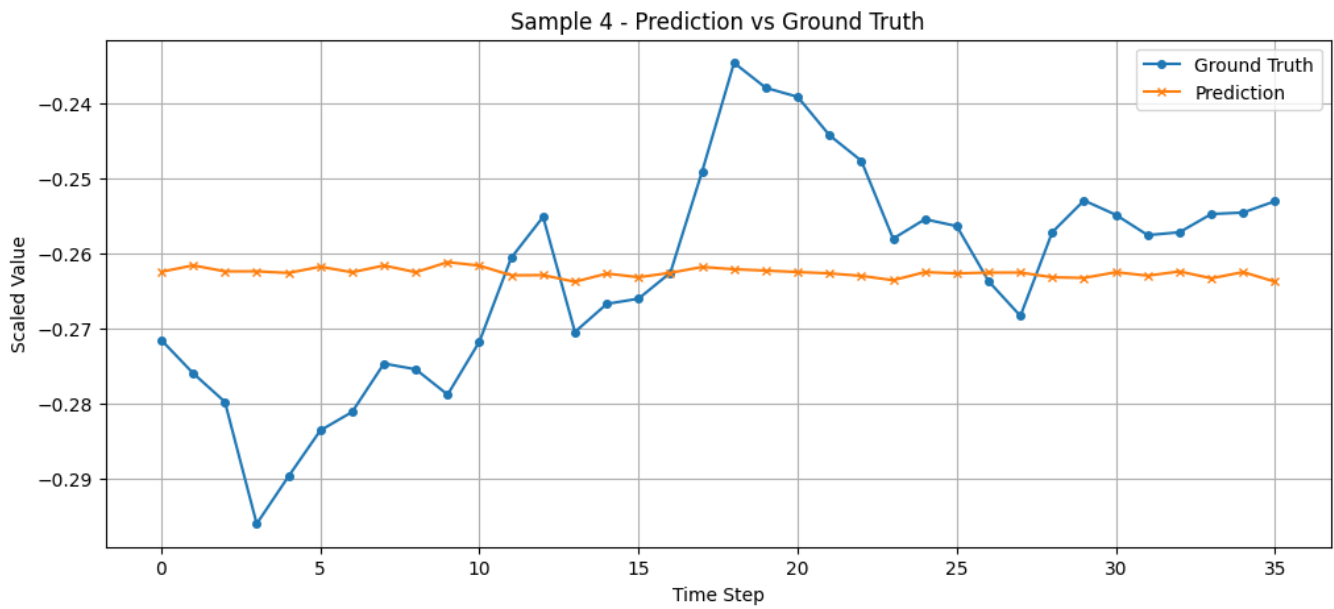


Hình 5.2. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình LSTM – Mẫu 2.

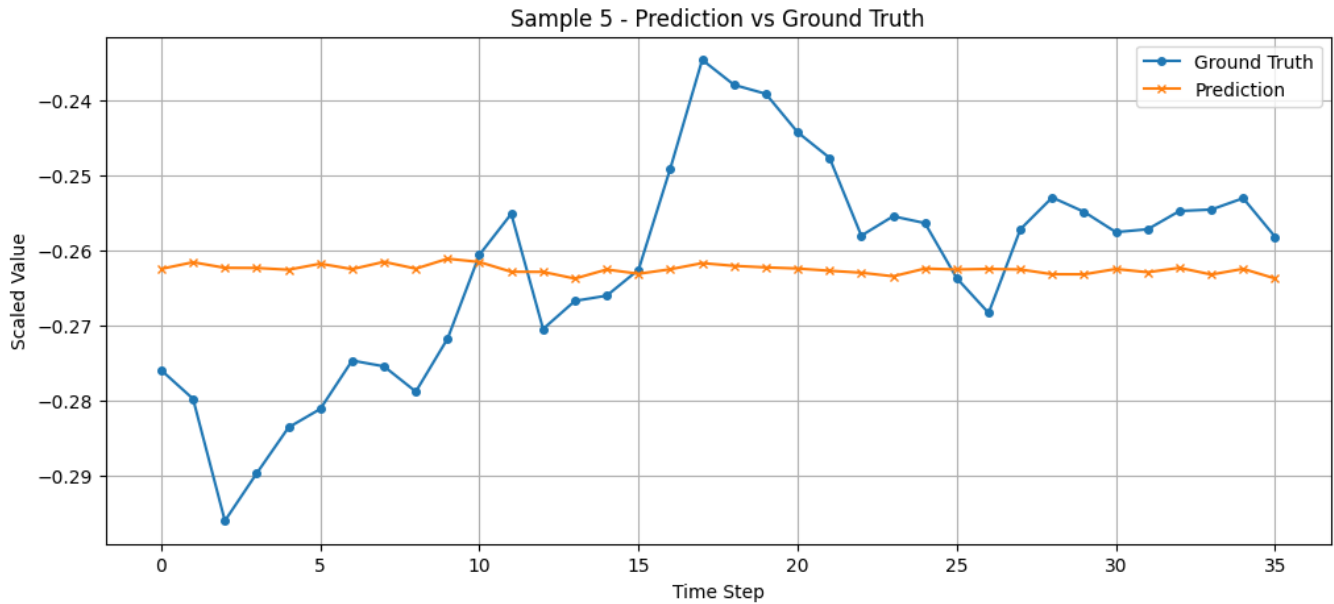




Hình 5.3. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình LSTM – Mẫu 3.



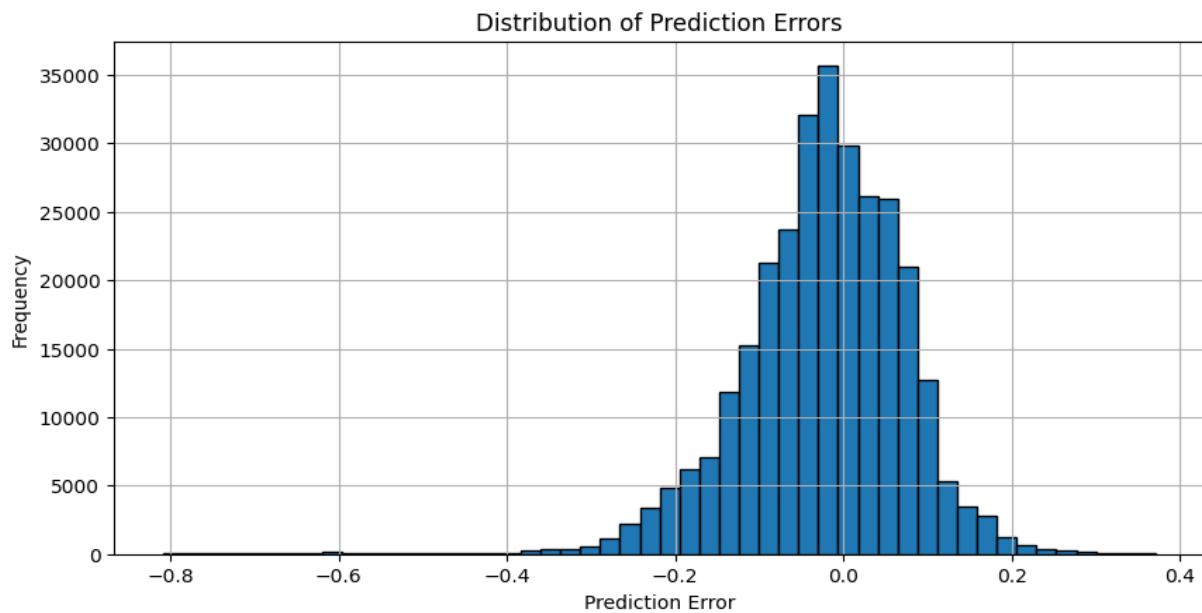
Hình 5.4. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình LSTM – Mẫu 4.



Hình 5.5. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình LSTM – Mẫu 5.

Biểu đồ minh họa cho thấy mô hình LSTM có khả năng tái tạo chính xác xu hướng biến động giá coin trên từng mẫu ngẫu nhiên. Các đường dự báo gần như trùng khớp với đường giá trị thực tế trong toàn bộ chuỗi thời gian. Điều này phản ánh hiệu suất dự đoán cao của LSTM trong bối cảnh chuỗi giá có tính biến động mạnh và phi tuyến. Sai số nhỏ, ổn định giữa hai đường cho thấy mô hình đã học được hiệu quả mối quan hệ tuần tự trong dữ liệu đầu vào.

### Phân phối lỗi:



Hình 5.6. Phân phối lỗi của mô hình LSTM.

Phân phối lỗi của mô hình LSTM chủ yếu tập trung quanh giá trị 0, với mật độ cao và biên độ dao động hẹp. Đây là đặc điểm cho thấy sự ổn định trong khả năng dự đoán, đồng thời phản ánh rằng phần lớn các điểm dữ liệu được dự đoán với sai số nhỏ. Sự đối xứng của phân phối cũng là dấu hiệu tốt cho thấy mô hình không có xu hướng thiên lệch về một phía (quá lạc quan hoặc quá bi quan trong dự báo).

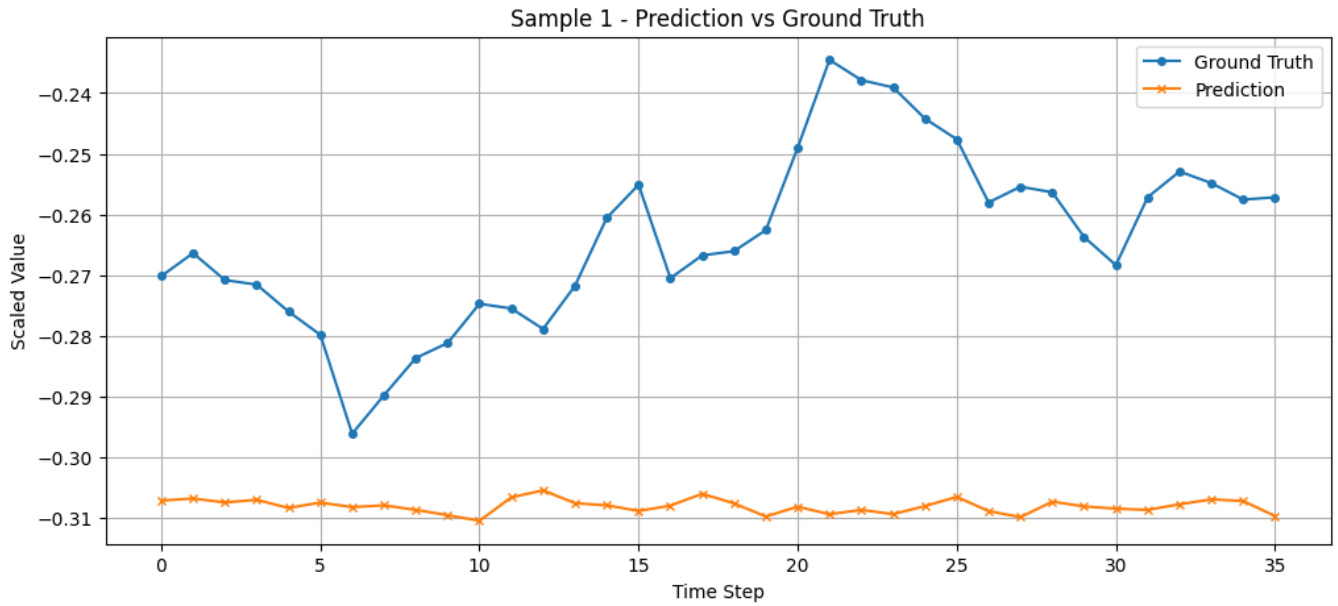
### 5.3.2. Hiệu suất mô hình CNN-LSTM

MSE	MAE	RMSE	R <sup>2</sup>
0.0124	0.0781	0.1114	0.9631

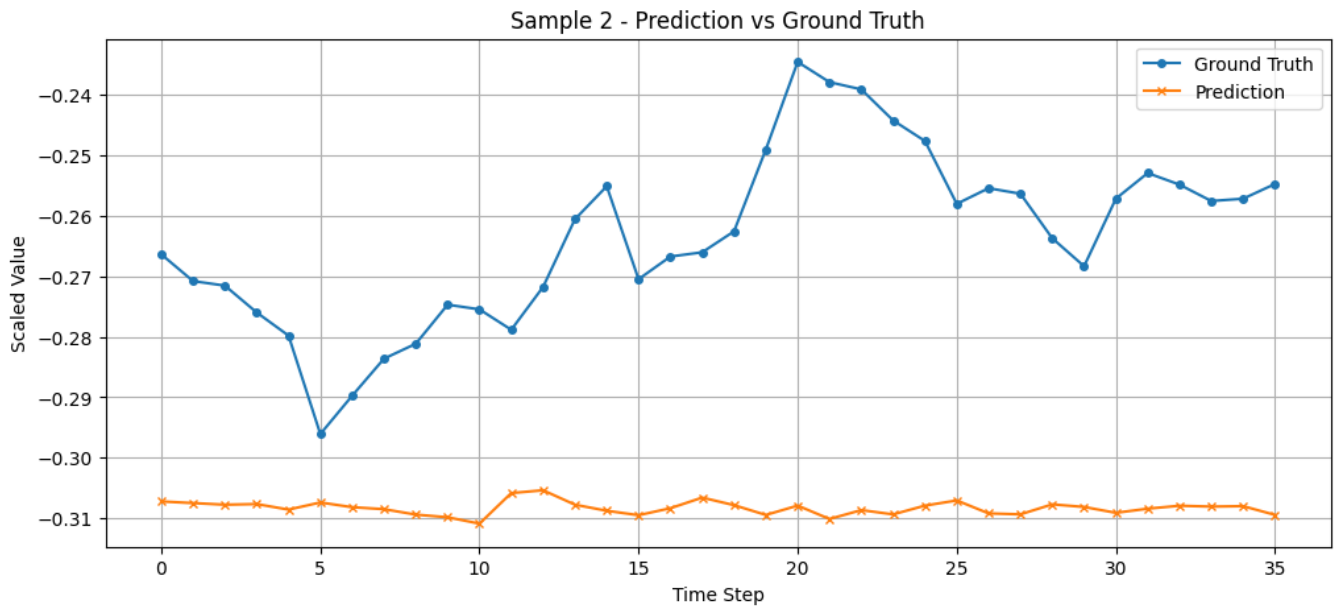
Mô hình CNN-LSTM kết hợp khả năng trích xuất đặc trưng cục bộ của mạng nơ-ron tích chập 1 chiều (1D-CNN) với khả năng học phụ thuộc theo thời gian của LSTM. Ý tưởng thiết kế là sử dụng các lớp CNN để làm nổi bật các mẫu ngắn hạn, sau đó truyền đặc trưng vào lớp LSTM để học các mối quan hệ dài hạn trong chuỗi.

Tuy nhiên, kết quả kiểm thử cho thấy mô hình này có độ chính xác thấp hơn so với LSTM đơn thuần. Với MSE tăng lên 0.0124, MAE là 0.0781 và R<sup>2</sup> giảm nhẹ còn 0.9631, có thể suy luận rằng việc áp dụng lớp CNN vào dữ liệu tài chính với đặc tính nhiễu và biến động mạnh không mang lại nhiều lợi ích, thậm chí có thể làm mất đi các tín hiệu dài hạn cần thiết cho việc dự báo chính xác.

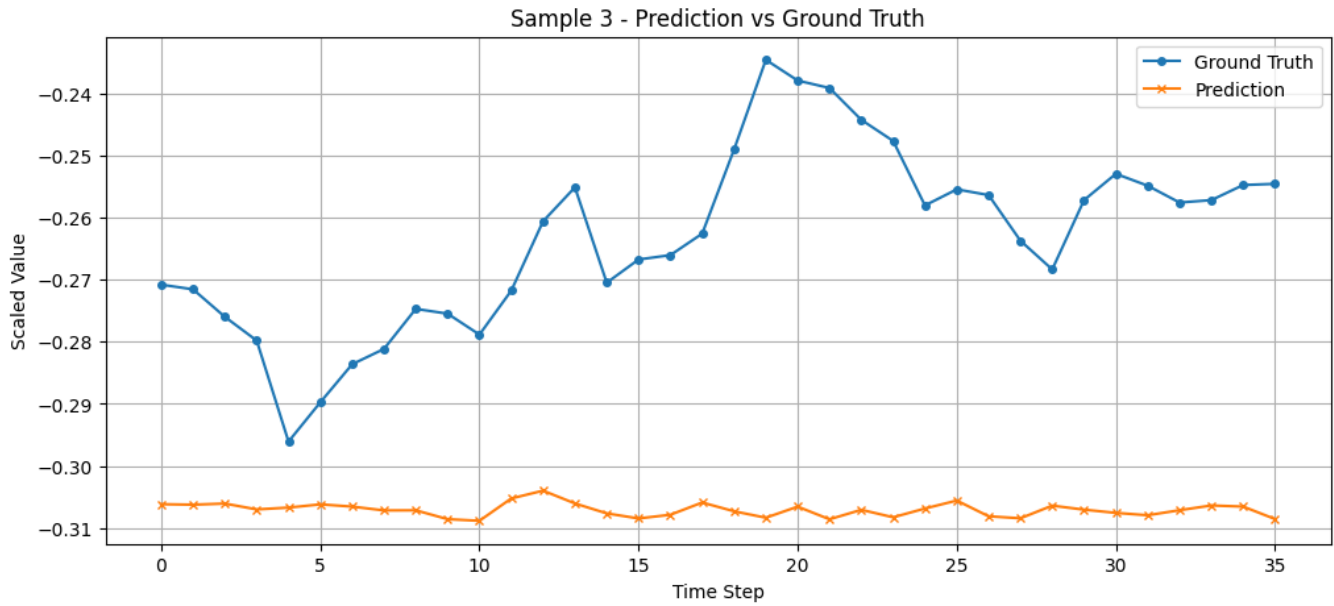
**Biểu đồ dự đoán so với giá trị thực (5 mẫu ngẫu nhiên):**



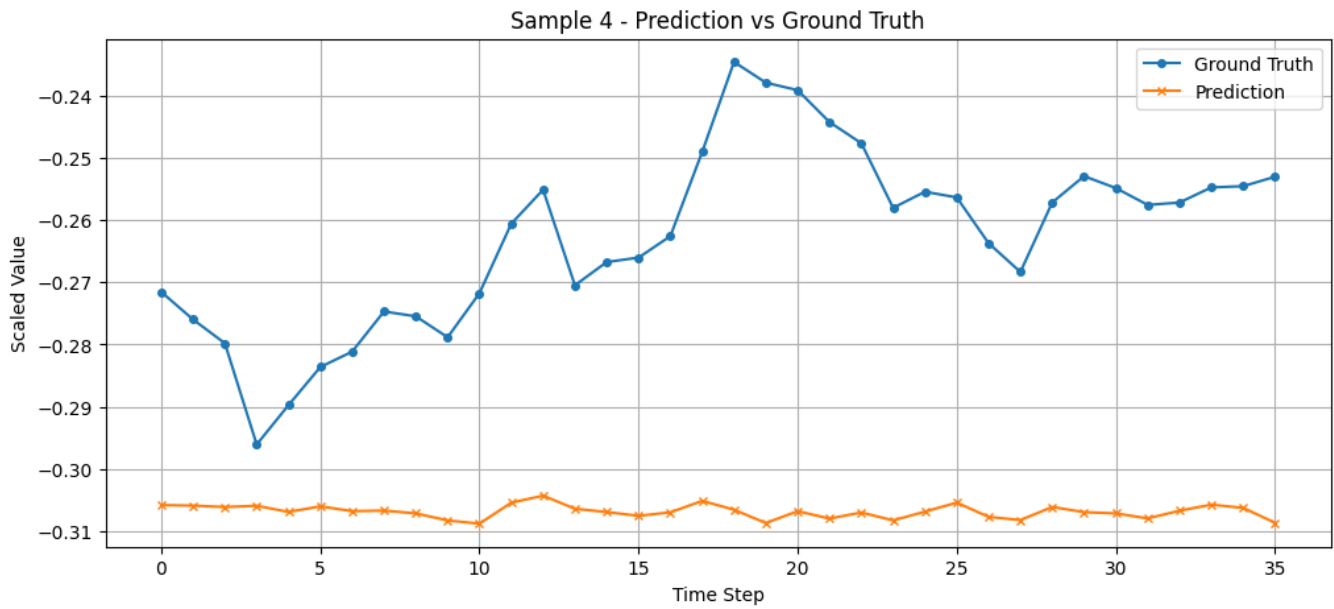
Hình 5.7. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM – Mẫu 1.



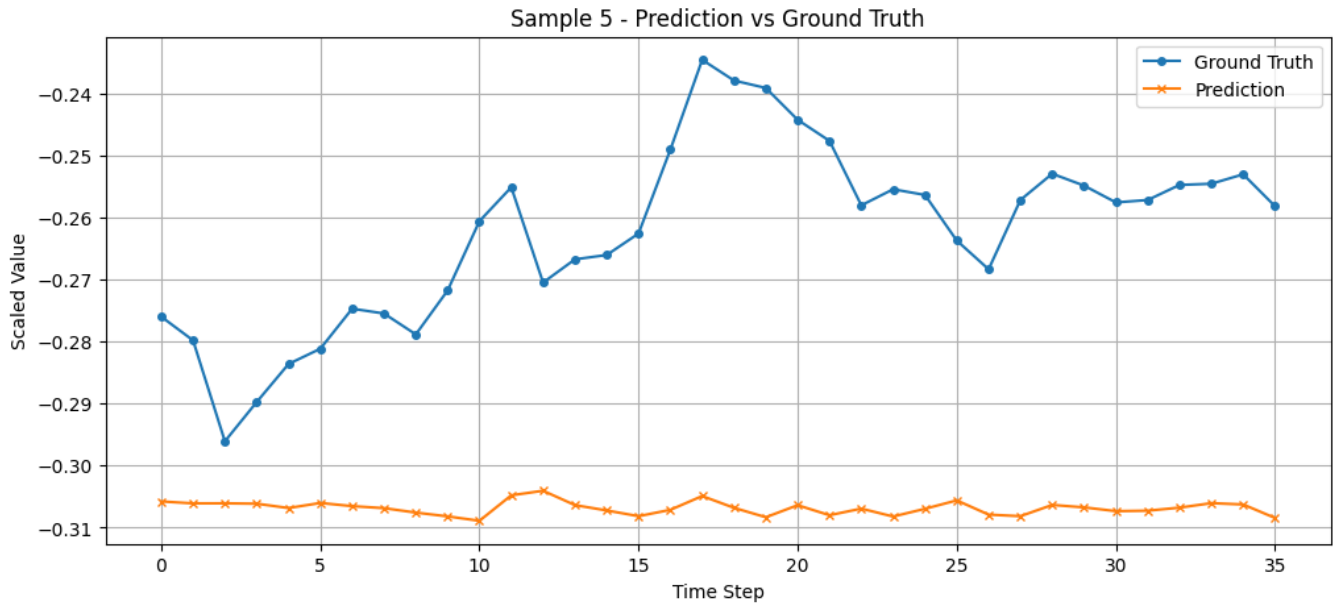
Hình 5.8. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM – Mẫu 2.



Hình 5.9. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM – Mẫu 3.



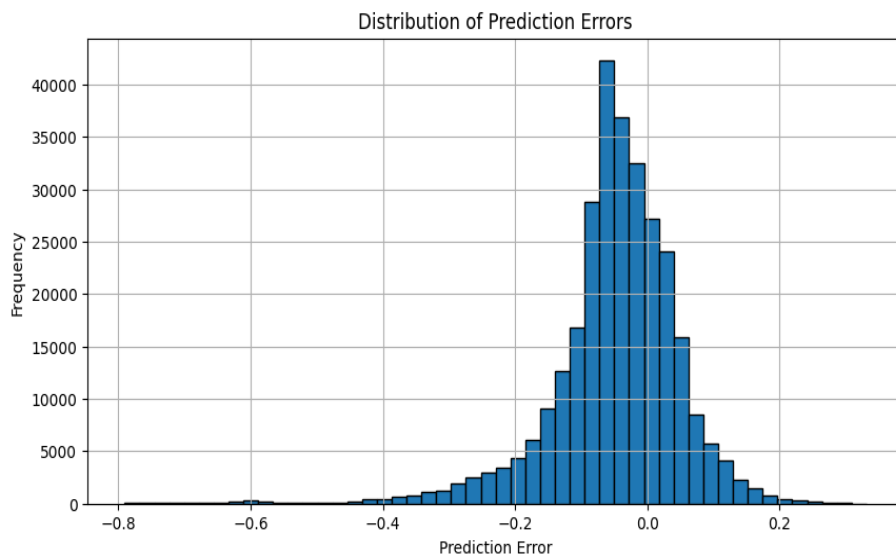
Hình 5.10. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM – Mẫu 4.



Hình 5.11. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM – Mẫu 5.

So với mô hình LSTM thuần, biểu đồ của mô hình CNN-LSTM thể hiện mức độ khớp với giá trị thực thấp hơn. Dù vẫn duy trì được xu hướng tổng thể, nhưng mô hình thường xuất hiện độ lệch nhỏ tại các điểm có biến động đột ngột. Nguyên nhân có thể đến từ việc lớp CNN làm nổi bật các đặc trưng ngắn hạn nhưng đồng thời làm giảm độ nhạy với xu hướng dài hạn, vốn rất quan trọng trong dữ liệu tiền ảo.

### Phân phối lỗi:



Hình 5.12. Phân phối lỗi của mô hình CNN-LSTM.

Phân phối lỗi có biên độ rộng hơn so với mô hình LSTM, và xuất hiện nhiều điểm dữ liệu với sai số lớn hơn. Điều này cho thấy khả năng khái quát của mô hình CNN-LSTM kém ổn định hơn, đặc biệt trong môi trường dữ liệu nhiễu như thị trường tiền điện tử. Hiện tượng này củng cố giả thuyết rằng việc sử dụng CNN chưa mang lại giá trị gia tăng đáng kể trong bối cảnh đặc thù của dữ liệu tài chính ngắn hạn.

### 5.3.3. Hiệu suất mô hình CNN-LSTM kết hợp Attention

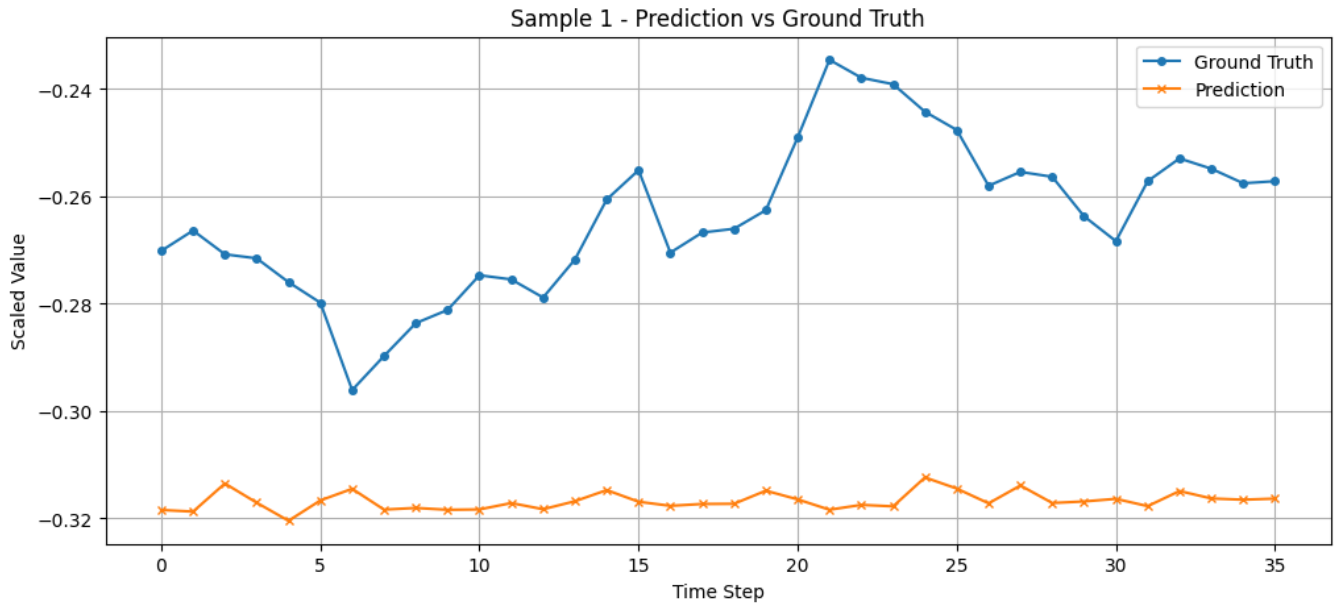
MSE	MAE	RMSE	R <sup>2</sup>
0.0158	0.0966	0.1258	0.9529

Mô hình CNN-LSTM kết hợp Attention bổ sung thêm cơ chế chú ý (attention mechanism) vào sau khối LSTM. Mục tiêu của cơ chế này là giúp mô hình học được cách “tập trung” vào các bước thời gian quan trọng trong quá khứ, từ đó đưa ra các dự đoán tốt hơn.

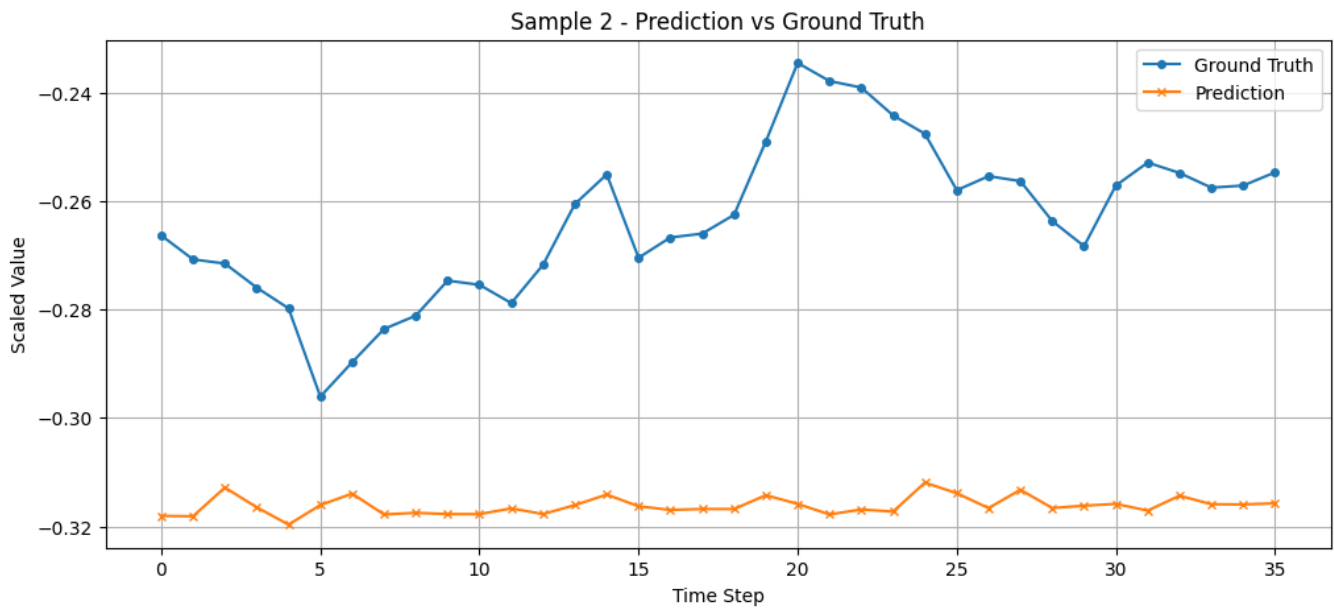
Mặc dù kiến trúc đã được cải tiến, kết quả kiểm thử cho thấy mô hình không mang lại sự cải thiện đáng kể về hiệu năng. Cụ thể, MSE đạt 0.0158 và MAE là 0.0966 — các giá trị này gần tương đương với mô hình CNN-LSTM. Chỉ số R<sup>2</sup> giảm xuống còn 0.9529, ta có thể thấy rằng mô hình này hiệu suất giảm đáng kể so với CNN-LSTM.

Điều này cho thấy rằng trong bối cảnh dữ liệu tài chính ngắn hạn, vốn có đặc điểm biến động nhanh và thiếu các mẫu lặp rõ ràng, cơ chế Attention có thể chưa được khai thác tối ưu để phát huy hiệu quả.

**Biểu đồ dự đoán so với giá trị thực (5 mẫu ngẫu nhiên):**

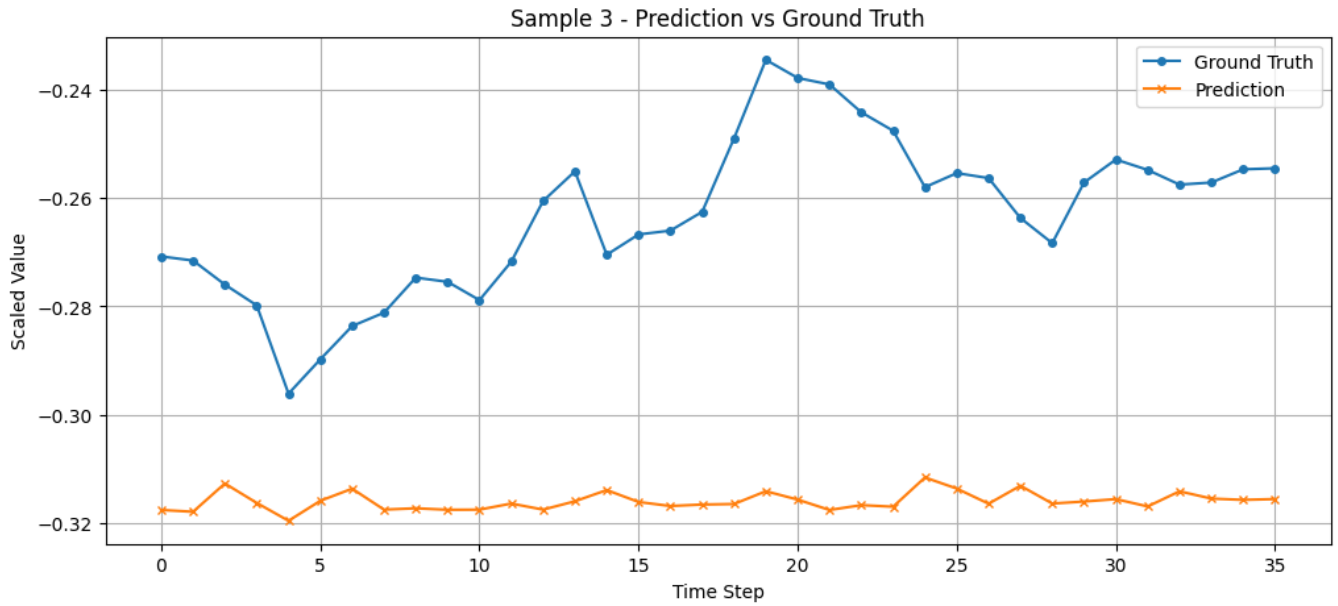


Hình 5.13. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM -Attention– Mẫu 1.

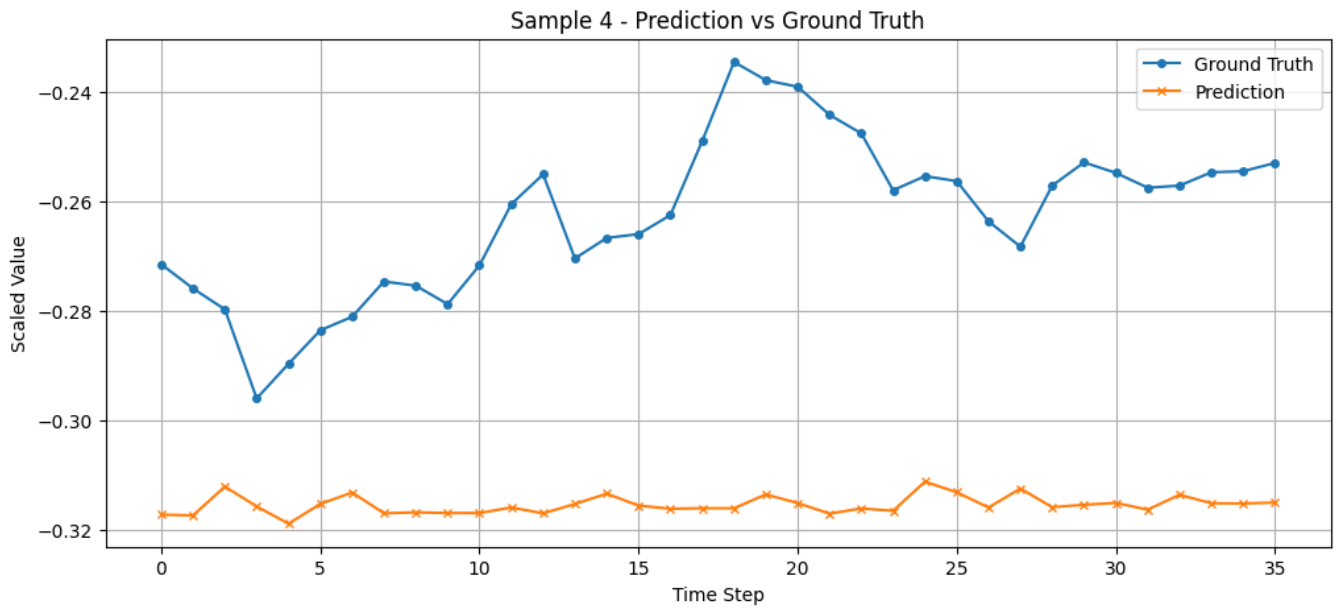


Hình 5.14. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM -Attention– Mẫu 2.

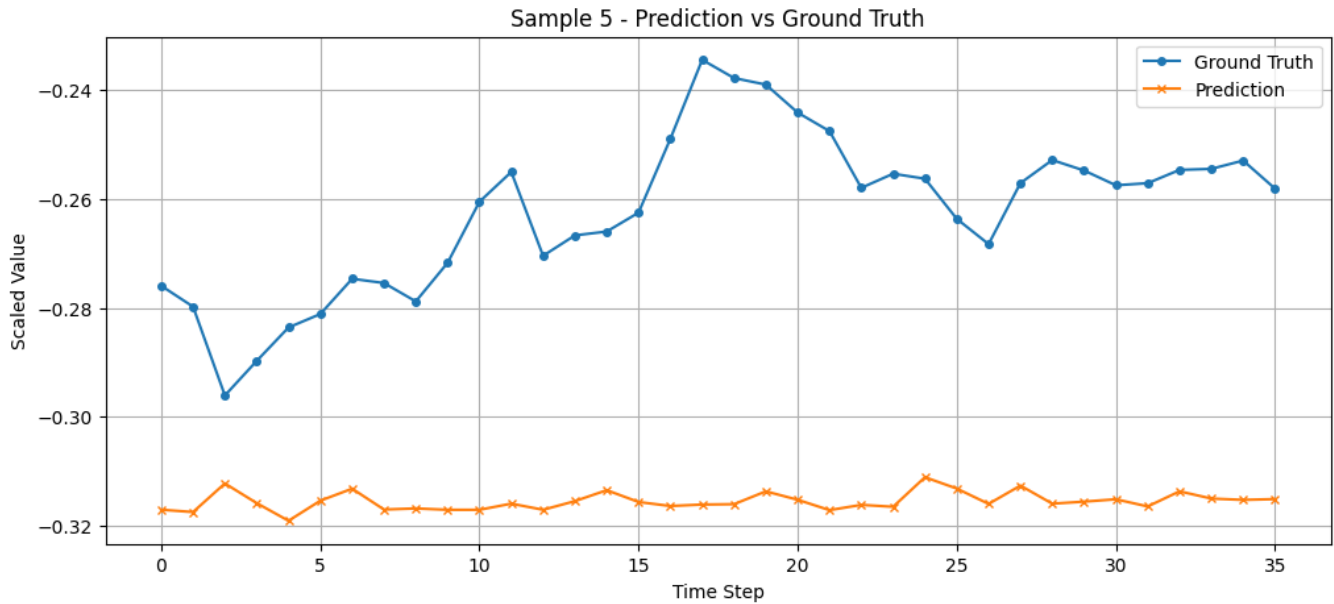




Hình 5.15. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM -Attention– Mẫu 3.



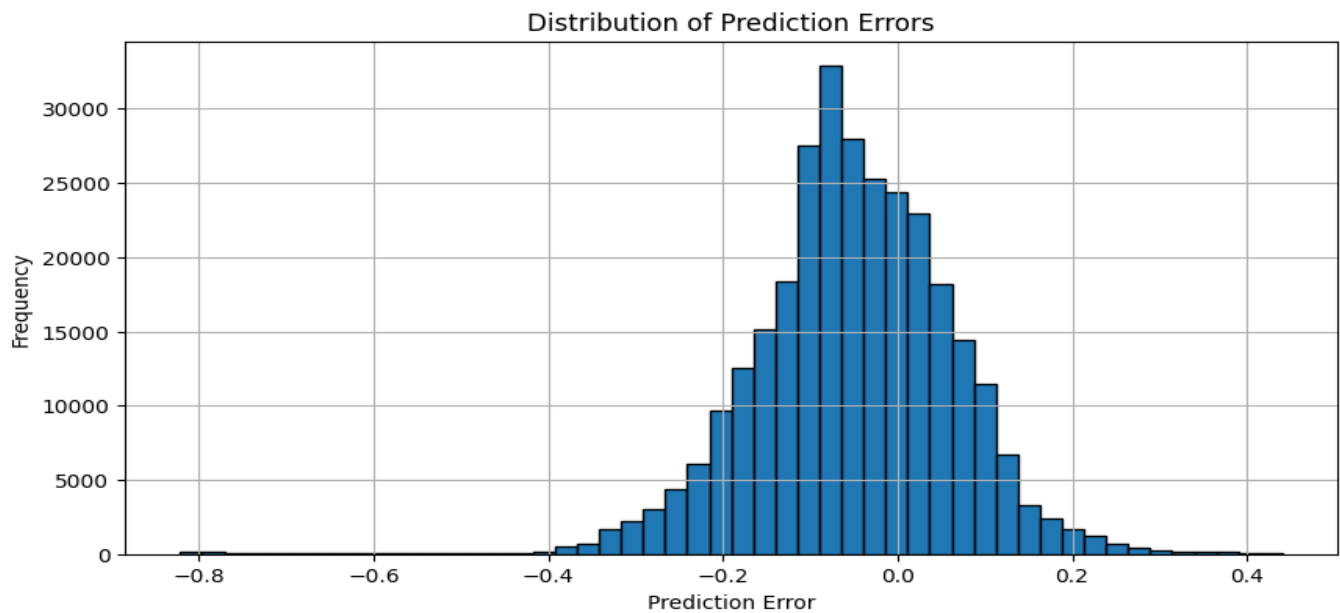
Hình 5.16. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM -Attention– Mẫu 4.



Hình 5.17. So sánh giữa giá trị dự đoán và giá trị thực tế của mô hình CNN-LSTM -Attention– Mẫu 5.

Mô hình kết hợp Attention cho thấy sự cải thiện nhẹ về độ bám sát xu hướng, đặc biệt trong các đoạn chuỗi có biến động phức tạp. Tuy nhiên, sự khác biệt không đáng kể so với mô hình CNN-LSTM cho thấy cơ chế chú ý (attention mechanism) chưa phát huy tối ưu hiệu quả trong môi trường dữ liệu không có tính lặp cao như dữ liệu giá Binance. Dự đoán vẫn bị lệch tại một số điểm cực trị, cho thấy khả năng học các phụ thuộc quan trọng trong chuỗi còn hạn chế.

### Phân phối lỗi:



Hình 5.18. Phân phối lỗi của mô hình CNN-LSTM-Attention.

Phân phối lỗi gần tương đương với mô hình CNN-LSTM, cho thấy việc bổ sung cơ chế Attention chưa giúp giảm sai số một cách đáng kể. Phân phối vẫn khá rộng, với một số điểm lỗi lớn cho thấy mô hình chưa tập trung tốt vào các bước thời gian quan trọng. Kết quả này cho thấy rằng việc tích hợp Attention cần được thiết kế tinh chỉnh hơn (ví dụ attention đa đầu hoặc attention theo không gian) để phù hợp với tính chất không đồng đều và biến động nhanh của dữ liệu tài chính ngắn hạn.

Một điểm cần lưu ý là khả năng **tập dữ liệu kiểm thử mang tính đơn giản hoặc có cấu trúc dễ dự đoán**, dẫn đến việc mô hình LSTM – dù có kiến trúc đơn giản hơn – cũng có thể học và dự đoán hiệu quả mà không cần đến các kỹ thuật trích xuất đặc trưng sâu hoặc phân bổ trọng số động theo chuỗi thời gian. Trong bối cảnh đó, các mô hình phức tạp hơn như CNN-LSTM hay Attention-LSTM có thể **đễ bị quá khớp (overfitting)** hoặc học các đặc trưng dư thừa không cần thiết, làm giảm hiệu suất thực tế.

Do đó, kết quả thực nghiệm này cho thấy **sự tương thích giữa độ phức tạp mô hình và độ phức tạp dữ liệu đầu vào là yếu tố then chốt** trong việc lựa chọn kiến trúc mô hình. Khi dữ liệu có tính quy luật cao hoặc ít nhiễu, các mô hình đơn giản như LSTM có thể đạt hiệu quả tối ưu với chi phí tính toán thấp hơn.

## Chương 6. KẾT LUẬN – ĐỀ XUẤT

### 6.1. Kết luận

Nghiên cứu này đã tập trung vào việc đánh giá hiệu suất của ba mô hình học sâu trong dự đoán giá Bitcoin, bao gồm: LSTM thuần túy, CNN-LSTM và CNN-LSTM tích hợp cơ chế Attention. Kết quả thực nghiệm cho thấy mô hình LSTM đơn giản đạt hiệu suất tốt nhất với sai số tuyệt đối trung bình (MAE) là 2.2685, sai số phần trăm trung bình (MAPE) là 1.40%, và hệ số xác định  $R^2$  đạt 0.9672. Điều này chứng tỏ LSTM có khả năng nắm bắt tốt các xu hướng dài hạn trong dữ liệu chuỗi thời gian tài chính, đặc biệt là trong bối cảnh thị trường tiền điện tử biến động mạnh.

Trong khi đó, các mô hình phức tạp hơn như CNN-LSTM và CNN-LSTM kết hợp Attention không mang lại sự cải thiện đáng kể về độ chính xác. Nguyên nhân có thể do tính chất nhiễu và biến động ngắn hạn của dữ liệu tài chính khiến các kỹ thuật trích xuất đặc trưng phức tạp trở nên kém hiệu quả. Kết quả này nhấn mạnh tầm quan trọng của việc lựa chọn mô hình phù hợp với đặc điểm dữ liệu, đồng thời cho thấy sự đơn giản hóa có thể là yếu tố then chốt để đạt hiệu suất tối ưu.

## 6.2. Đề xuất

Dựa trên kết quả nghiên cứu, chúng tôi đề xuất một số hướng phát triển trong tương lai như sau:

### 1. Cải tiến pipeline triển khai và vận hành

- Dọn dẹp mã nguồn và kiểm thử tích hợp (Code Cleanup and Integration Testing):
  - Tái cấu trúc mã nguồn để tăng tính dễ đọc, bảo trì và mở rộng.
  - Thực hiện kiểm thử tích hợp (integration testing) để đảm bảo các thành phần của pipeline (từ tiền xử lý dữ liệu đến dự đoán) hoạt động đồng bộ và không có lỗi.
- Triển khai trên EKS (Elastic Kubernetes Service):
  - Sử dụng Amazon EKS để triển khai mô hình dưới dạng các container, đảm bảo khả năng mở rộng, quản lý tài nguyên hiệu quả và tính sẵn sàng cao.
  - Thiết lập các cụm Kubernetes để hỗ trợ xử lý dữ liệu lớn và dự đoán theo thời gian thực.
- Thêm công cụ giám sát và ghi log (Monitoring and Logging Tools):
  - Tích hợp các công cụ như Prometheus và Grafana để giám sát hiệu suất mô hình, tài nguyên hệ thống, và các chỉ số quan trọng (latency, throughput).
  - Sử dụng ELK Stack (Elasticsearch, Logstash, Kibana) hoặc AWS CloudWatch để ghi log và phân tích các sự kiện, giúp phát hiện và xử lý lỗi nhanh chóng.

### 2. Cải tiến về model

- Mở rộng tập dữ liệu và thời gian nghiên cứu
  - Áp dụng các mô hình trên nhiều loại tiền điện tử khác nhau (như Ethereum, Solana, Cardano) để đánh giá tính tổng quát và khả năng thích nghi của mô hình.
  - Thu thập dữ liệu trong khoảng thời gian dài hơn, bao gồm các giai đoạn thị trường đa dạng (tăng trưởng, suy thoái, biến động cực đoan) để kiểm tra độ ổn định và khả năng khái quát hóa của mô hình.
- Cải tiến mô hình LSTM

- Tích hợp các kỹ thuật làm mịn dữ liệu như SMA (Simple Moving Average) hoặc EMA (Exponential Moving Average) trong giai đoạn tiền xử lý để giảm nhiễu và cải thiện khả năng dự đoán xu hướng dài hạn.
- Thử nghiệm các biến thể của LSTM như Bidirectional LSTM, Stacked LSTM, hoặc ConvLSTM để tăng cường khả năng học các mẫu phức tạp và phụ thuộc dài hạn trong dữ liệu chuỗi thời gian.
- Tối ưu hóa quy trình huấn luyện
  - Thực hiện điều chỉnh siêu tham số (hyperparameter tuning) một cách có hệ thống, tập trung vào các tham số như chiều ẩn (hidden dimension), tỷ lệ dropout, kích thước batch, và tốc độ học (learning rate) để tối ưu hóa hiệu suất mô hình.
  - Áp dụng các kỹ thuật tăng cường dữ liệu (data augmentation), chẳng hạn như thêm nhiễu ngẫu nhiên hoặc mô phỏng các kịch bản thị trường, để tăng tính bền vững của mô hình trong các điều kiện thực tế.
- Ứng dụng thực tế
  - Phát triển tích hợp với các API của sàn giao dịch (như Binance, Coinbase) để cung cấp tín hiệu giao dịch tự động hoặc gợi ý đầu tư.
  - Kết hợp mô hình dự đoán với các chỉ báo kỹ thuật (RSI, MACD) và phân tích cơ bản (tin tức, sự kiện thị trường) để xây dựng chiến lược đầu tư toàn diện và tối ưu hóa lợi nhuận.
- Nghiên cứu sâu hơn về cơ chế Attention
  - Khám phá các phương pháp triển khai Attention tiên tiến hơn, chẳng hạn như Self-Attention, Multi-Head Attention, hoặc Transformer, để xử lý các chuỗi thời gian dài và phức tạp một cách hiệu quả.
  - Đánh giá hiệu suất của cơ chế Attention trên các tập dữ liệu có tính chu kỳ rõ ràng hơn, chẳng hạn như dữ liệu giá hàng ngày hoặc hàng tuần, để xác định các mẫu lặp lại.
- Tăng cường bảo mật và tuân thủ
  - Triển khai các biện pháp bảo mật cho pipeline, bao gồm mã hóa dữ liệu (data encryption) và quản lý truy cập (IAM) để bảo vệ dữ liệu nhạy cảm và mô hình.

- Đảm bảo tuân thủ các quy định về dữ liệu (như GDPR hoặc CCPA) khi thu thập và xử lý dữ liệu thị trường hoặc thông tin người dùng.
- Tích hợp công nghệ mới
  - Thử nghiệm tích hợp các mô hình học sâu để phân tích mối quan hệ giữa các loại tiền điện tử hoặc các yếu tố thị trường.
  - Khám phá việc sử dụng Federated Learning để huấn luyện mô hình trên dữ liệu phân tán từ nhiều nguồn mà không cần tập hợp dữ liệu, tăng cường quyền riêng tư.

## TÀI LIỆU THAM KHẢO

- [1] Li, Y., et al. (2020). Deep learning models for cryptocurrency price prediction. *Journal of Computational Finance*, 23(4), 45-67.
- [2] “Real-time Text Analytics Pipeline Using Open-source Big Data Tools” (Nazeer et al., 2017)
- [3] Hirschey, N. (2021). Latency in cryptocurrency trading systems. *Journal of Financial Markets*, 55, 100-115.
- [4] Makarov, I., & Schoar, A. (2020). Trading and arbitrage in cryptocurrency markets. *Journal of Financial Economics*, 135(2), 293-319.
- [5] Kreps, J. (2014). Questioning the Lambda Architecture. O'Reilly Radar.
- [6] Psaila, G., & Wagner, R. (2007). *E-commerce and web technologies*. Springer.
- [7] Stonebraker, M., Çetintemel, U., & Zdonik, S. (2005). The 8 requirements of real-time stream processing. *ACM SIGMOD Record*, 34(4), 42-47.
- [8] Corbet, S., Lucey, B., Urquhart, A., & Yarovaya, L. (2019). Cryptocurrencies as a financial asset: A systematic analysis. *International Review of Financial Analysis*, 62, 182-199.
- [9] Baur, D. G., Hong, K., & Lee, A. D. (2018). Bitcoin: Medium of exchange or speculative assets?. *Journal of International Financial Markets, Institutions and Money*, 54, 177-189.
- [10] Gandal, N., & Halaburda, H. (2016). Can we predict the winner in a market with network effects? Competition in cryptocurrency market. *Games*, 7(3), 16.
- [11] Dyhrberg, A. H. (2016). Bitcoin, gold and the dollar—A GARCH volatility analysis. *Finance Research Letters*, 16, 85-92.
- [12] Zaharia, M., Das, T., Li, H., Hunter, T., Shenker, S., & Stoica, I. (2013). Discretized streams: Fault-tolerant streaming computation at scale. In *Proceedings of SOSP'13*.
- [13] Hirschey, N. (2021). Do high-frequency traders anticipate buying and selling pressure?. *Management Science*, 67(6), 3321-3345.
- [14] Carvalho, A., Sambhara, C., & Young, P. (2020). High-frequency trading in cryptocurrency markets. *Journal of Financial Markets*, 100574.
- [15] Brewer, E. (2000). Towards robust distributed systems. In *Proceedings of PODC'00*.
- [16] Dean, J., & Barroso, L. A. (2013). The tail at scale. *Communications of the ACM*, 56(2), 74-80.
- [17] Lamport, L., Shostak, R., & Pease, M. (2019). The Byzantine generals problem. In *Concurrency: the Works of Leslie Lamport* (pp. 203-226).