

# Exploring Diffusion Models for Generative Forecasting of Financial Charts

Taeyeon Lee<sup>1,2</sup>, Jiwon Park<sup>2</sup>, Kyunga Bang<sup>2</sup>, Seunghyun Hwang<sup>1,2,3</sup>, Ung-Jin Jang<sup>1\*</sup>

<sup>1</sup>FnGuide Inc.

<sup>2</sup>GenAI in Finance, MODULABS

<sup>3</sup>Department of Applied Data Science, Sungkyunkwan University

taeyeonglee@fnguide.com, mary000605@ewha.ac.kr,  
{kbang1002, hsh1030}@g.skku.edu, coorung77@fnguide.com

## Abstract

Recent advances in generative models have enabled significant progress in tasks such as generating and editing images from text, as well as creating videos from text prompts, and these methods are being applied across various fields. However, in the financial domain, there may still be a reliance on time-series data and a continued focus on transformer models, rather than on diverse applications of generative models. In this paper, we propose a novel approach that leverages text-to-image model by treating time-series data as a single image pattern, thereby enabling the prediction of stock price trends. Unlike prior methods that focus on learning and classifying chart patterns using architectures such as ResNet or ViT, we experiment with generating the next chart image from the current chart image and an instruction prompt using diffusion models. Furthermore, we introduce a simple method for evaluating the generated chart image against ground truth image. We highlight the potential of leveraging text-to-image generative models in the financial domain, and our findings motivate further research to address the current limitations and expand their applicability.

## 1 Introduction

Recent advances in Large Language Models (LLMs)[4, 1, 18] and diffusion-based generative models [3, 17, 8, 9] have enabled powerful cross-modal capabilities such as generating and editing images from text. In finance, however, most stock price prediction [12, 21] still relies on single-modality time-series models like LSTMs [6] or Transformers [19], which may struggle to capture visually discernible signals (chart patterns, candlestick shapes) [20, 11, 22] and to integrate heterogeneous factors such as news or sentiment. Existing image-based studies [20, 15, 2] typically perform classification with architectures like ResNet [7] or ViT [5], without explicitly modeling the temporal evolution of chart patterns and relying only on chart snapshots without integrating signals such as volume or technical indicators.

In real markets, human traders often interpret price movements through visual patterns—head-and-shoulders, double bottoms, candlestick wicks—because these configurations implicitly encode market psychology such as fear, greed, and indecision [13, 14]. Traditional time-series approaches, while effective at modeling sequential dependencies, lose the structural “shape” context of price movements, struggle to jointly represent multiple heterogeneous signals (e.g., price, volume, indicators) [10], and offer limited interpretability from a trader’s perspective.

---

\*Corresponding author

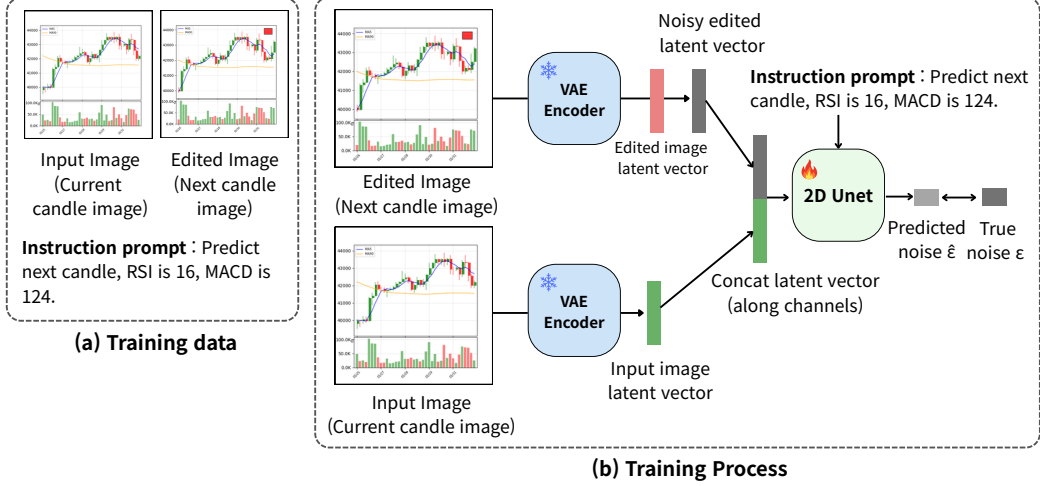


Figure 1: **Overall architecture of our method.** (a) Training data: We construct a paired dataset consisting of an input image, an edited image, and an instruction prompt. (b) Training Process : We encode the input image and the edited image, and then fine-tune the 2D U-Net of Stable Diffusion with the instruction prompt.

In this work, we explore a novel approach that generates the next chart image from the current chart and an instruction prompt using a text-to-image diffusion model [17]. To support this, we construct paired datasets of input charts, instruction prompts, and edited chart images, and fine-tune a U-Net in latent space. Because the generated outputs are stochastic and differ from time-series predictions, we further introduce a simple image-marking evaluation that compares generated images to ground-truth charts via RGB analysis.

We conducted preliminary experiments on cryptocurrency data, and the results suggest that our method is feasible and promising, though performance remains limited. Our method represents an exploratory step toward leveraging generative models in financial forecasting, and suggests potential extensions such as incorporating sentiment signals from news articles.

Our main contributions are as follows:

- We propose a novel approach that leverages a text-to-image generative model, treating time-series financial data as visual patterns for next-step chart generation.
- We introduce a paired dataset construction method combining chart images, instruction prompts, and edited images, enabling the model to learn both technical indicators (RSI, MACD) and visual chart patterns.
- We introduce a simple image-marking evaluation method to quantify prediction accuracy from generated images. Also we show experiment results on cryptocurrency data, and discuss the limitations and opportunities for extending this approach to richer multi-modal signals such as sentiment.

## 2 Method

### 2.1 Training Dataset

Inspired by InstructPix2Pix [3], we propose a simple yet effective method for fine-tuning Stable Diffusion [17] to learn chart patterns and generate next-candle charts based on given instructions. As shown in Figure 1 (a), we construct paired datasets consisting of an input image, an edited image, and an instruction prompt. In our dataset, the input image is represented as a 4-hour candlestick chart at the  $n$ -th timestep, incorporating trading volume as well as the SMA5 and SMA90 lines. This allows the model to learn not only from the candlestick chart itself, but also from the relationships with trading volume and moving averages. The edited image is constructed as the candlestick chart at timestep  $n+3$ , with an evaluation marker placed at the upper-right corner of the image. A red mark is assigned if the price increases by more than 2%, a blue mark if it decreases, and a black mark

Table 1: Category-wise classification performance. We generate our images for the test input image, and then evaluate our model by analyzing the RGB values to determine its consistency with the ground truth.

Category	Precision (%)	Recall (%)	F1-Score (%)	Support
blue	8.06	11.63	9.52	43
red	9.26	14.93	11.43	67
black	85.60	77.94	81.59	671
Overall Acc.	68.89% (538/781)			

Table 2: Confusion matrix of classification results. P. blue, P. red, and P. black indicate the predicted classes, while A. blue, A. red, and A. black indicate the actual classes.

	P. blue	P. red	P. black
A. blue	5	4	34
A. red	3	10	54
A. black	54	94	523

otherwise. The instruction prompt consists of the RSI and MACD values at timestep  $n$ , formatted with the prefix: `Predict next candle`.

## 2.2 Training Process

As illustrated in Figure 1 (b), we freeze the VAE of Stable Diffusion 1.5 [17] and fine-tune only the 2D U-Net. Given an input image  $\mathbf{I}_{in}$  (at timestep  $n$ ) and an edited image  $\mathbf{I}_{edit}$  (at timestep  $n+3$ ), the frozen VAE encoder produces latent vectors  $\mathbf{l}_z, \mathbf{n}_e \in \mathbb{R}^{4 \times 64 \times 64}$ , where  $\mathbf{n}_e$  is further perturbed with gaussian noise. The two latents are concatenated along the channel dimension to form  $\mathbf{x} \in \mathbb{R}^{8 \times 64 \times 64}$ , and the U-Net input stem is modified accordingly.

The instruction prompt provides the technical indicators at timestep  $n$ , formatted as: `Predict next candle, RSI is {value}, MACD is {value}`, and is injected via cross-attention. The U-Net is trained to predict the noise  $\hat{\epsilon}_\theta$  from the noisy latent  $\mathbf{x}_t$  and instruction  $\mathbf{c}$ , with the standard denoising objective:

$$\mathcal{L}(\theta) = \mathbb{E}_{\mathbf{x}, \mathbf{c}, t, \epsilon} [\|\epsilon - \text{UNet}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2].$$

Through this training process, our method leverages the image generation capability of Stable Diffusion [17] to learn chart-specific patterns while incorporating instruction prompts, thereby enabling the generation of future chart images from the current chart.

## 2.3 Inference for Next Chart Generation

Since we fine-tune Stable Diffusion 1.5 [17], the model preserves its image generation capability while learning chart-specific patterns. Given the current chart image with RSI and MACD values as an instruction prompt, the model can generate a candlestick chart four hours ahead. The generated chart includes the evaluation mark, trading volume, and moving averages, similar to the edited image, as illustrated in Figure 2.

# 3 Experiments

## 3.1 Experimental Setup

**Dataset.** Our method represents time series data as image charts, making it possible to fine-tune for any asset. We focus particularly on Bitcoin price, as its data is easy to collect and its patterns are clearly visible. Following Section 2.1, we construct a paired dataset consisting of 2,419 training data pairs in total. We use 4-hour candlestick charts of Bitcoin future prices with Binance exchange from January 1, 2024, to March 1, 2025, as the input images. Each training image contains a total of 40 candlesticks. The edited image is constructed in the same manner as the input image, but with an evaluation marker placed at the upper-right corner.

**Implementation detail.** We fine-tune Stable Diffusion 1.5 [17] using an NVIDIA A100 80GB GPU, with a batch size of 16, a gradient accumulation step of 4, and a total of 28,000 optimization steps. Our inference step is set to 20, image guidance scale to 1.0, and guidance scale to 2.0. We use the EulerAncestralDiscreteScheduler as the image scheduler.

**Evaluation.** We construct an evaluation dataset using 4-hour candlestick charts of Bitcoin futures from the Binance exchange, covering the period from March 17, 2025 to July 31, 2025. Our evaluation

(a) **Instruction prompt** : Predict next candle, RSI is 83.68, MACD is 534.52.

(b) **Instruction prompt** : Predict next candle, RSI is 20.6, MACD is 217.8.

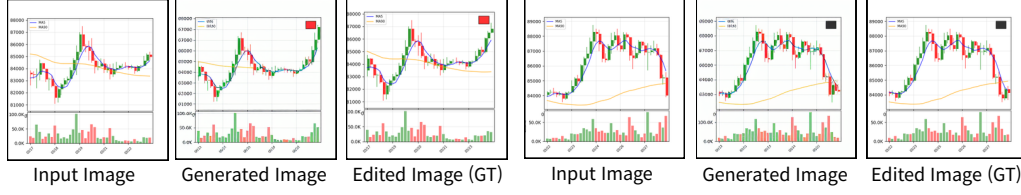


Figure 2: **Generated samples.** We generate the image from an input image and an instruction prompt, while the edited image serves as the ground truth. Our results demonstrate that the model is capable of sufficiently learning the visual patterns of charts.

dataset consists of three paired components: input image, instruction prompt, and edited image, for a total of 781 pairs. We simply classify the generated image by reading the RGB values of this mark region and mapping them to one of the three classes via a color-thresholding rule. Accuracy is computed as the fraction of samples where the predicted class matches the ground-truth class.

### 3.2 Quantitative Results

Table 1 presents the classification results, with an overall accuracy of 68.89% (538/781). The performance is mainly driven by the black class (F1-score 81.59%), while the blue and red classes show much lower F1-scores (9.52% and 11.43%), as confirmed by the confusion matrix in Table 2. Many blue and red samples are misclassified as black, indicating that minority-class patterns are less distinctive and easily absorbed into the majority category.

These findings suggest that the proposed generative approach may capture broad trend patterns, while it appears less effective in distinguishing finer variations across underrepresented classes. Contributing factors include dataset imbalance, the difficulty of diffusion models in retaining chart-specific structures, and imperfect alignment between generated and ground truth images.

Although the results leave room for improvement, our method is intended as an exploratory step. Our experiments show that text-to-image generative models can encode meaningful chart-like structures, and we believe that methods such as class balancing, domain-aware conditioning, or hybrid generative–discriminative frameworks could improve quantitative performance in future work.

### 3.3 Qualitative Results

As shown in Figure 2, we use the input image and instruction prompt as inputs. Based on these, our model generates the next candlestick chart, which can be seen in the generated image. The generated image shares similarities with the edited image, and we can observe that it is generated by taking into account trading volume, moving averages, as well as RSI and MACD indicators. Despite its limited quantitative performance, the model can generate chart patterns and diverse samples to visualize and anticipate scenarios, opening possibilities for generative forecasting charts. We expect traders could use it as a scenario simulation tool to visually explore market psychology.

## 4 Limitation and Conclusion

**Limitation.** Our approach still suffers from limitation. The RGB-based evaluation remains simplistic and may overlook the true financial validity of generated charts. The dataset size is limited and centered on cryptocurrency, restricting broader applicability. In addition, the predictive performance of our model is still modest, highlighting that current generative approaches are not yet competitive with traditional forecasting methods. Nevertheless, we believe that future work holds significant opportunities for improvement. For example, instruction prompts could be enriched with diverse external signals, such as financial news, FOMC announcements, or market sentiment reports. Multi chart inputs (e.g., combining 15-minute and 4-hour candlesticks) could provide the model with richer temporal context. More powerful generative backbones, such as Stable Diffusion XL [16], may improve both fidelity and consistency. Beyond this, domain-specific evaluation metrics, larger and more diverse datasets, and multimodal integration with order book depth or social media sentiment could substantially enhance both accuracy and interpretability.

In conclusion, while our current model shows limited predictive performance, this work illustrates a novel direction: reframing time-series forecasting in finance through the lens of generative models. By treating financial data as visual patterns for chart generation, we open up a research pathway that blends generative modeling with technical and multimodal analysis, paving the way for more advanced and practically useful financial forecasting systems.

## Acknowledgments

This research was supported by Brian Impact Foundation, a non-profit organization dedicated to the advancement of science and technology for all.

## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Jeongseok Bang and Doojin Ryu. Cnn-based stock price forecasting by stock chart images. *Romanian Journal of Economic Forecasting*, 26(3):120–128, 2023.
- [3] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. *arXiv preprint arXiv:2211.09800*, 2022.
- [4] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198, 2024.
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [6] Alex Graves. Long short-term memory. *Supervised sequence labelling with recurrent neural networks*, pages 37–45, 2012.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [8] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022.
- [9] Taegyeong Lee, Jeonghun Kang, Hyeonyu Kim, and Taehwan Kim. Generating realistic images from in-the-wild sounds. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7160–7170, 2023.
- [10] Mengxia Liang, Shaocong Wu, Xiaolong Wang, and Qingcai Chen. A stock time series forecasting approach incorporating candlestick patterns and sequence similarity. *Expert Systems with Applications*, 205:117595, 2022.
- [11] Yaohu Lin, Shancun Liu, Haijun Yang, Harris Wu, and Bingbing Jiang. Improving stock trading decisions based on pattern recognition using machine learning technology. *PloS one*, 16(8):e0255558, 2021.
- [12] Siqi Lu, Buyao Song, and Guowen Li. Enhancing multi-factor stock selection with transformer networks: A comparative analysis against traditional machine learning models. *Procedia Computer Science*, 266:1028–1034, 2025.
- [13] Filip Martinsson and Ivan Liljeqvist. Short-term stock market prediction based on candlestick pattern analysis, 2017.
- [14] Edrees Ramadan Mersal, Kürşat Mustafa Karaoğlu, and Hakan Kutucu. Enhancing market trend prediction using convolutional neural networks on japanese candlestick patterns. *PeerJ Computer Science*, 11:e2719, 2025.
- [15] Jakub Pizoń, Łukasz Kański, Jan Chadam, and Bartłomiej Pęk. Image-based time series trend classification using deep learning: A candlestick chart approach. *Advances in Science and Technology. Research Journal*, 19(11), 2025.

- [16] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- [17] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- [18] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [19] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [20] Marc Velay and Fabrice Daniel. Stock chart pattern recognition with deep learning. *arXiv preprint arXiv:1808.00418*, 2018.
- [21] Jue Xiao, Tingting Deng, and Shuochen Bi. Comparative analysis of lstm, gru, and transformer models for stock price prediction. In *proceedings of the international conference on digital economy, blockchain and artificial intelligence*, pages 103–108, 2024.
- [22] Wenke Zhu, Weisi Dai, Chunling Tang, Guoxiong Zhou, Zewei Liu, and Yunjing Zhao. Pmanet: a time series forecasting model for chinese stock price prediction. *Scientific Reports*, 14(1):18351, 2024.