

A Data Science Pipeline for Algorithmic Trading: A Comparative Study of Applications for Finance and Cryptoeconomics

1st Luyao Zhang

Data Science Research Center and Social Science Division
Duke Kunshan University
Suzhou, China
lz183@duke.edu

1st Tianyu Wu

Duke Kunshan University
Suzhou, China

1st Saad Lahrichi

Duke Kunshan University
Suzhou, China

2nd Carlos-Gustavo Salas-Flores

Duke Kunshan University
Suzhou, China

2nd Jiayi Li

Duke Kunshan University
Suzhou, China

Abstract—Recent advances in Artificial Intelligence (AI) have made algorithmic trading play a central role in finance. However, current research and applications are disconnected information islands. We propose a generally applicable pipeline for designing, programming, and evaluating the algorithmic trading of stock and crypto assets. Moreover, we demonstrate how our data science pipeline works with respect to four conventional algorithms: the moving average crossover, volume-weighted average price, sentiment analysis, and statistical arbitrage algorithms. Our study offers a systematic way to program, evaluate, and compare different trading strategies. Furthermore, we implement our algorithms through object-oriented programming in Python3, which serves as open-source software for future academic research and applications.

Index Terms—algorithmic trading, data science pipeline, finance, cryptoeconomics, open-source software, Python

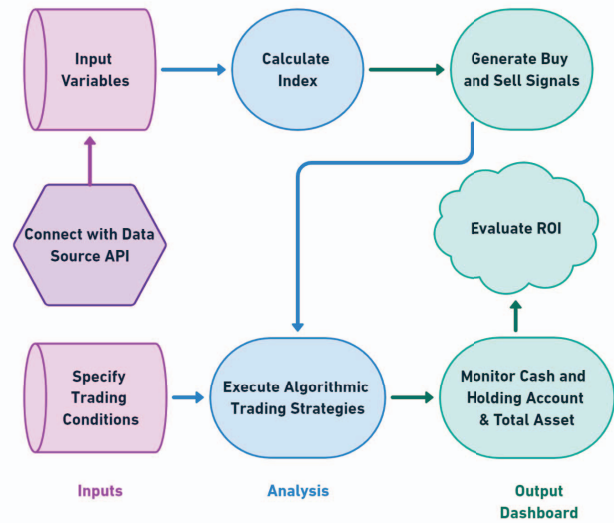
I. INTRODUCTION

Recent advances in AI have made algorithmic trading (AT) play a central role in finance: Burgess [1] estimates that AT accounted for 80% of the entire equity turnover in the U.S. by transaction volume in 2018, up from 50% in 2011. AT empowers traders with cutting-edge technological innovation so that they can execute trades with predefined strategies, capturing profitable opportunities faster and with a low attention cost [2]. Moreover, AT contributes to market efficiency and liquidity [3]. However, current research and applications are disconnected information islands [4]. AT in financial industries is generally a black box for which scientific evaluation is almost impossible [5]. Academic research has more clearly elaborated on algorithms. However, current research diverges in many ways. The lack of process consistency makes ceteris paribus comparison difficult. Moreover, there is little open-source software for designing AT strategies, which leads to

Tianyu Wu, Saad Lahrichi, Jiayi Li, and Carlos-Gustavo Salas-Flores were supported by the Summer Research Scholar Program at Duke Kunshan University as research affiliates in Prof. Luyao Zhang's project entitled "How Fintech Empowers Asset Valuation: Theory and Applications"

unnecessary obstacles to collaborative learning, research, and innovation.

Fig. 1. The data science pipeline for algorithmic trading: The general workflow includes inputs (pink), analysis (blue), and the output dashboard (green). At the input stage, we first connect with the data source API to input variables for calculating indices that are necessary to calculate buy-and-sell signals and then specify trading conditions. At the analysis stage, we calculate indices that are necessary to calculate buy-and-sell signals and execute algorithmic trading strategies. At the output stage, we visualize three dashboards: (1) the time series of buy-and-sell signals, (2) the cash and holding accounts and total assets, and (3) the return on investment (ROI).



We propose a generally applicable pipeline for designing, programming, and evaluating the algorithmic trading of stock and crypto tokens. Figure 1 represents the general workflow: inputs, analysis, and the output dashboard. We first connect with the data source API to input variables for calculating indices that are necessary to calculate buy-and-sell signals; then,

we further specify trading conditions to execute algorithmic trading strategies based on the signals. Finally, we generate visualizations to monitor cash and holding accounts and evaluate the return on investment (ROI). We demonstrate that our data science pipeline is generally applicable to conventional algorithms, including the moving average crossover, volume-weighted average price, sentiment analysis, and statistical arbitrage algorithms.

Our study offers a systematic way to program and compare different trading strategies. Furthermore, we implement our algorithms by object-oriented programming in Python3, which serves as open-source software for future academic research and applications. We introduce the methodology applied to two algorithms in Section II, present the data and results in Section III, and discuss the results in Section IV.

II. METHODOLOGY

A. Moving Average Crossover

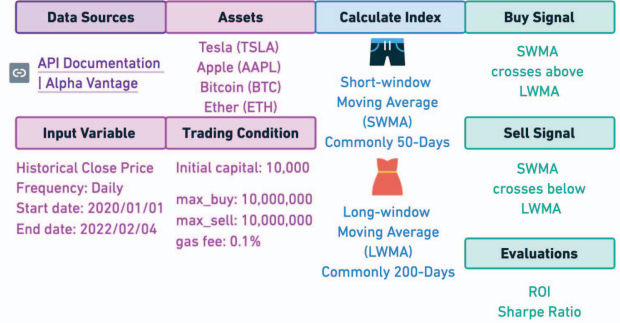
The simple moving average (SMA), first proposed by Joseph E. Granville Granville [6], is an indicator of the trend of stock prices [7, 8]. As in Equation 1, the SMA with window n is the average of the stock price in the past n days.

$$SMA_t^n = \frac{1}{n} \sum_{i=t-n+1}^t p_i \quad (1)$$

A moving average crossover occurs when a short-window moving average (SWMA) and a long-window moving average (LWMA) intersect. On the one hand, the SWMA is an indicator of recent market prices. On the other hand, the LWMA represents the long-term or equilibrium price. A buy signal appears when the SWMA (such as the 50-day moving average) crosses above the LWMA (such as the 200-day moving average), predicting a bull market, and this is where the golden crossover occurs. The embedded idea is that as long-term indicators carry more weight, the golden crossover indicates a bull market on the horizon and is reinforced by high trading volumes. Similarly, a sell signal is released when the SWMA crosses below the LWMA, predicting a bear market, and this is where the death crossover happens. Some researchers have already tested this trading rule on some commonly used indices in the US stock market, finding that this technical analysis can generate a higher return on investment than the buy-and-hold strategy; however, it still cannot reflect a downturn in a timely manner due to its lag property [9].

Figure 2 represents the general pipeline for the moving average crossover strategy. We first input historical closing prices through the Alpha Vantage API, from which we can obtain all daily closing price data from the US stock market from the time to market, but we select 40% of the historical data here because we would like to make a comparison based on the return on investment (RoI) and Sharpe ratio more intuitively in the Discussion section. Then, we calculate the SWMA and LWMA for generating buy-and-sell signals. Next, we simulate the strategies with an initial capital of 10,000 USD

Fig. 2. **Moving average crossover:** The algorithm starts from the left with inputs (pink), from which we obtain our data from the Alpha Vantage API (purple). The raw input variables are historical daily closing prices. Then, we perform our moving average analysis (blue) to determine the buy-or-sell signals (green). Lastly, we evaluate the performance by using the ROI and Sharpe ratio as metrics.



and set a maximum capacity to enable buying and selling at each crossover. Here, for simplicity, we ignore the transaction fee. Finally, we evaluate the performance of this strategy by comparing the ROI and Sharpe ratio to a buy-and-hold strategy.

B. Volume-Weighted Average Price

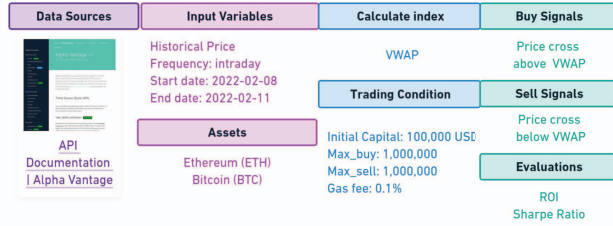
The SMA provides the recent price trend of a security. However, it does not take into account the level of volume traded at each price level. On the other hand, the volume-weighted average price (VWAP) gives the average price for intraday trading weighted by the transaction volume at each price. Equation (3) represents the formal definition of the VWAP, where P_t is the price at time t and Q_t is the corresponding volume traded at that price.

$$VWAP = \frac{\sum P_t \cdot Q_t}{\sum Q_t} \quad (2)$$

Cesari et al. [10] suggest that traders have incorporated institutional investors' tendency into target the VWAP into their intraday trading strategies [10]. A bullish/bearish market is indicated when prices cross above/below the VWAP and traders shall buy/sell accordingly [11, 12, 13]. Menkhoff (2010) indicates that the performance of VWAP strategies varies significantly and that the VWAP does not necessarily outperform the buy-and-hold strategy [14].

Figure 3 represents the pipeline for a VWAP case study. We keep most parts consistent with the pipeline for the SMA. Unfortunately, the VWAP index is directly available from

Fig. 3. **Summary of the VWAP trading algorithm:** The data sources used are shown in purple. We relied on the Alpha Vantage API to collect the necessary data. The input variables used and the assets we trade are in pink. The trading conditions and the calculated index (in our case, the VWAP) are shown in blue. Green represents the conditions for buy/sell signals as well as the metrics used for evaluating the algorithm.



the Alpha Vantage API¹ only for traditional stocks. For cryptocurrencies, we perform the calculation ourselves.

III. DATA AND RESULTS

A. Data

1) *Moving Average Crossover Data:* Table I shows all the data used in the moving average crossover algorithm:

TABLE I
MOVING AVERAGE CROSSOVER: DATA

DATA SOURCE: ALPHA VANTAGE API. THIS TABLE SHOWS THE RAW DATA THAT WE RETRIEVED FROM THE API. THE PRICE OF STOCKS AND CRYPTO ASSETS CORRESPONDS TO THE CLOSING PRICE OF THE GIVEN ASSET ON A GIVEN DATE.

Variable	Frequency	Unit	Description
Date	daily	YYYY-MM-DD	Date and time for which the data were recorded
Close	daily	USD	Price at which the stock ended trading in a given time period
Short MA	daily	USD	Average price of a security within a certain period, typically 50 days.
Long MA	daily	USD	Average price of a security within a certain period, typically 200 days.
Signal	-	-	Buy-and-sell signal (e.g., TSLA, AAPL for stock, BTC, ETH for crypto)

2) *Volume-Weighted Average Price:* Using the Alpha Vantage API, we collected intraday stock data derived from securities information processor (SIP) market-aggregated data

¹Alpha Vantage is a website that offers financial market data through its developer-friendly APIs. Alpha Vantage provides both traditional assets data and forex and cryptocurrency data. Upon installation, we used its time series stocks API to retrieve intraday time series of the stocks we explored. We also used the tech indicators API, from which we retrieved daily VWAP data for the same stocks. Without that API, we would have had to calculate the VWAP by applying its formula to the data. We have used Alpha Vantage because its API is among one of the few that offers free intraday data. Unlike the open-high-low-closing data that many websites offer, intraday data are rarer and usually are not free of charge. Quandl, for example, has high-quality intraday data, but they are not free. A key limitation of the use of Alpha Vantage's free API is that it only allows us to query only the last 15 days. It is worth noting that the documentation says that the TIME_SERIES_INTRADAY method can retrieve 1-2 months of intraday data, while the TIME_SERIES_INTRADAY_EXTENDED returns the trailing 2 years. As another important limitation, the maximum number of requests is 5 API requests per minute and 500 API requests per day. Alpha Vantage offers a premium API key, which allows for a higher call limit as well as more historical data.

as well as intraday cryptocurrency data. The interval chosen was 5 min. Table II summarizes the data. It shows the variables we use, the frequency at which these variables are collected, the unit of each variable, and a short description.

TABLE II
VOLUME WEIGHTED AVERAGE PRICE: DATA
DATA SOURCE: ALPHA VANTAGE API

Variable	Frequency	Unit	Description
Date	5 min	YYYY-MM-DD HH:MM:SS	Date and time for which the data were recorded
Close	5 min	USD	Price at which the stock ended trading in a given time period
VWAP	5 min	USD	Average price of a security within a day, adjusted for its volume. Available for an API call only for traditional stocks; manually calculated using the formula for crypto.
Ticker	-	-	Stock symbol (e.g., TSLA, AAPL for traditional, BTC, ETH for crypto)
Interval	5 min	min/hr/day	Time difference between two data points

B. Results

1) *Moving Average Crossover:* Here, we present one typical cryptocurrency, Ether (ETH), to generate buy-and-sell signals, followed by the moving average crossover strategies. We visualize our holdings and cash flow during this testing period, and we determine and compare the performance of the moving average crossover to the simple buy-and-hold strategy.

Figure 4 is a visualization of how the buy-and-sell strategy is defined when applying the moving average strategy to perform the backtesting on ETH. According to Figure 2, when the short-window moving average (SWMA) crosses above the long-window moving average (LWMA), it generates a buy signal. Conversely, it produces a sell signal. The signal is specifically drawn on the closing price line.

Fig. 4. **Buy-and-Sell Signal: ETH moving average crossover**



Figure 5 shows how the portfolio would change in the time series after applying the crossover strategy to claim to buy or sell ETH during the backtesting period. Here, we assume that the transaction fee is 0.1% per transaction. We carry out buying behavior by buying the maximum number of shares with all the cash held, and we carry out selling behavior by selling all the shares to obtain cash.

From Figure 6, it is obvious that over the two-year backtesting period, the total revenue increases.

From Figure 6, the ROI given this backtesting period is 849.8%, performing much better than the simple buy-and-hold strategy, which produces an ROI of 12.0%.

Fig. 5. Portfolio time series: ETH moving average crossover



Fig. 6. Gross ROI: ETH moving average crossover vs. buy-and-hold
Moving average crossover strategy ROI: 849.84%
Buy & hold strategy ROI: 11.97%



From Figure 7, the Sharpe ratio given this backtesting period is 1.69, showing that this strategy, which produces a Sharpe Ratio result of 3.24, has a higher ability to deal with financial risks.

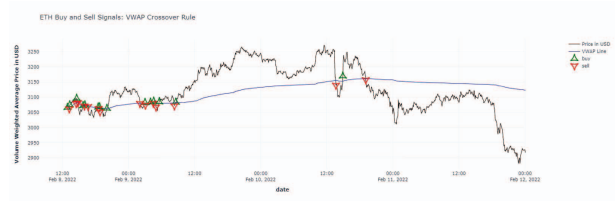
Fig. 7. Sharpe Ratio: ETH moving average crossover vs. buy-and-hold
Moving average crossover strategy Sharpe ratio: 0.98
Buy-and-hold strategy Sharpe ratio: 2.60



2) *Volume-Weighted Average Price*: Figure 8 shows the buy-and-sell signals generated using the VWAP strategy when trading ETH. The black line shows the evolution of the price of ETH for the period retrieved. The blue line shows the VWAP value over the same period. The value of the VWAP was calculated for each 5-min interval from the price and closing data. We see that there is an upward green arrow signifying a buy signal each time the price of ETH crosses the VWAP line from below. We have downward red arrows for sell signals, which are generated whenever the price of ETH crosses the VWAP line from above. We see that when the current price and the VWAP line intersect repeatedly in a short time period, there are multiple consecutive buy-and-sell signals.

Figure 9 displays the evolution of the portfolio when trading ETH using the VWAP strategy. We plotted the current amount of cash in red, the value of the holdings in green, and the sum of both cash and holdings in red (total). As expected from the multiple buy-and-sell signals in Figure 8 early on, there are many fluctuations in the amount of cash and holdings on the

Fig. 8. Buy-and-Sell Signal: Volume-Weighted Average Price



left side of the graph. There are far fewer as we move to the right (as there are fewer signals produced there). The overall total seems to have decreased from its initial value of 100,000.

Fig. 9. Portfolio time series: Volume-Weighted Average Price



Figure 10 displays a comparison of the ROI when trading ETH using the buy-and-hold strategy (in black) and the VWAP strategy (in blue). The buy-and-hold strategy looks the same as the price evolution graph in Figure 8. For the studied time period, the trader would have an ROI of approximately -4%. Using the VWAP also results in a negative ROI, very close to the buy-and-sell ROI. For this time period, using the buy-and-sell strategy or the VWAP strategy would have resulted in a loss.

Fig. 10. Gross ROI: Volume-Weighted Average Price vs. Buy-and-Hold
Volume-weighted average price strategy ROI: -3.93%
Buy-and-hold strategy ROI: -4.26%

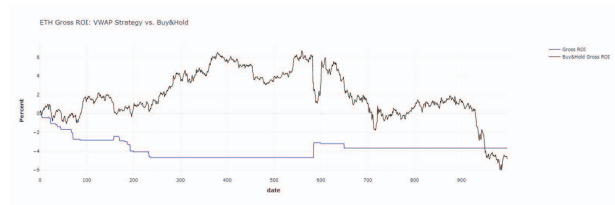


Figure 11 displays a comparison of the Sharpe ratio when trading ETH using the buy-and-hold strategy (in black) and the VWAP strategy (in blue). We see that both Sharpe ratios end up being negative, with the buy-and-hold ratio being slightly below zero, and the VWAP ratio being close to -2%. Both the buy-and-hold and VWAP strategies return negative Sharpe ratios in this case. We therefore cannot say that using the VWAP improves our returns, especially in the very short term.

IV. DISCUSSION

In conclusion, we propose a data science pipeline for algorithmic trading and verify that the pipeline is generally

Fig. 11. **Sharpe Ratio: Volume-Weighted Average Price vs Buy-and-Hold**
Volume-Weighted Average strategy Sharpe ratio: -1.62
Buy-and-hold strategy Sharpe ratio: -0.52



applicable for designing, programming, and evaluating the algorithmic trading of stock and crypto assets. Furthermore, we implement our algorithms through object-oriented programming in Python3, which serves as open-source software for future academic research and applications. The data, code, and supplementary results can be found on the Github repository: <https://github.com/SciEcon/SRS2021>.

We evaluated our data science pipeline based on two conventional algorithmic trading strategies.

- **The moving average crossover** is among the most prominent investment strategies in finance that build mathematical models to further forecast future market movements through historical market records [15]. Currently, researchers and industry practitioners are keenly interested in analyzing the moving average crossover strategy in crypto trading. Fang et al. [16] applied the simple moving average trading strategy based on daily closing price data for the 11 most traded cryptocurrencies over the period 2016-2018. Their results showed that technical trading rules, excluding Bitcoin, produced a desirable annualized excess return of 8.76%. Moreover, due to the high volatility of crypto prices, a revised approach to the moving average has started to be introduced by considering the tolerance of time ranges for the intersection of transaction signals in real industry practices [17].
- **The volume-weighted average price** can be viewed as an advanced version of the *moving average crossover strategy* that averages by weights of transaction volumes. Although the prior literature has mainly focused on using the VWAP as one of the technical indicators used to trade assets, few efforts have been made to use the index for cryptocurrencies. One emerging trend, however, is the adoption of the VWAP in the industry practices of DeFi applications. A prominent example is the Ampleforth DeFi protocol [18], which issues stable coin AMPL. Ampleforth uses VWAP data from decentralized Chainlink oracles [19] to expand or contract the supply of AMPL accordingly, every 24 hours, adjusting the price of AMPL to be approximately one US dollar regardless of the volatility in the market demand for AMPL [20].

We also provide an online appendix in the Github repository: <https://github.com/SciEcon/SRS2021>. In the online appendix, we present additional results of the aforementioned

two algorithms, the moving average crossover and volume-weighted average price, and we provide an additional evaluation of our data science pipeline based on two additional algorithms: **sentiment analysis** and **pairs trading**.

Our research is seminal in inspiring future innovations in two ways.

- First, the current research shows how our data science pipeline is generally applicable for implementing and evaluating existing trading algorithms. Furthermore, our approach can be applied to design, implement, and evaluate new trading algorithms, as in Liu and Zhang [21]. Since there currently does not exist a consensus on the valuation of crypto assets, effective algorithmic trading strategies for various crypto tokens have yet to be found.
- Second, our data science pipeline is general enough for researchers to customize various settings, such as data sources, input variables, assets, trading conditions, indices, buy-and-sell signals, and evaluation indicators. However, the pipeline might better serve more advanced reinforcement learning algorithms and sentiment analysis by incorporating existing general-purpose or domain-specific data science frameworks, as in [22].

ACKNOWLEDGMENTS

We thank the Alpha Vantage API for providing academic accounts for the data querying in our research. We thank Zesen Zhuang for his assistance in hosting tutorial sessions for Tianyu Wu, Saad Lahrichi, Carlos-Gustavo Salas-Flores, and Jiayi Li about object-oriented programming in Python3.

REFERENCES

- [1] N. Burgess, "An introduction to algorithmic trading: Opportunities & challenges within the systematic trading industry," *SSRN Electronic Journal*, 2019.
- [2] G. P. M. Virgilio, "High-frequency trading: a literature review," *Financial Markets and Portfolio Management*, vol. 33, pp. 183–208, 06 2019. [Online]. Available: <https://link.springer.com/article/10.1007/s11408-019-00331-6>
- [3] P. Gomber and K. Zimmermann, *Algorithmic Trading in Practice*, S.-H. Chen, M. Kaboudan, and Y.-R. Du, Eds. Oxford University Press, 02 2018.
- [4] N. Reznik and L. Pankratova, "High-frequency trade as a component of algorithmic trading: market consequences." [Online]. Available: http://ceur-ws.org/Vol-2104/paper_174.pdf
- [5] A. Azzutti, W.-G. Ringe, and H. S. Stiehl, "Machine learning, market manipulation and collusion on capital markets: Why the 'black box' matters," papers.ssrn.com, 02 2021. [Online]. Available: <https://ssrn.com/abstract=3788872>
- [6] J. E. Granville, *New key to stock market profits*. Prentice-Hall, 1963.
- [7] A. RAUDYS and Z. PABARSKAITE, "Optimising the smoothness and accuracy of moving average for stock price data," *Technological and Economic Development of Economy*, vol. 24, pp. 984–1003, 05 2018.

- [8] I. Bhattacharjee and P. Bhattacharja, "Stock price prediction: A comparative study between traditional statistical approach and machine learning approach," *2019 4th International Conference on Electrical Information and Communication Technology (EICT)*, 12 2019.
- [9] W. BROCK, J. LAKONISHOK, and B. LeBARON, "Simple technical trading rules and the stochastic properties of stock returns," *The Journal of Finance*, vol. 47, pp. 1731–1764, 12 1992. [Online]. Available: <https://onlinelibrary.wiley.com/doi/full/10.1111/j.1540-6261.1992.tb04681.x>
- [10] R. Cesari, M. Marzo, and P. Zagaglia, "Effective trade execution," *SSRN Electronic Journal*, 2012.
- [11] S. A. BERKOWITZ, D. E. LOGUE, and E. A. NOSER, "The total cost of transactions on the nyse," *The Journal of Finance*, vol. 43, pp. 97–112, 03 1988.
- [12] A. R. Admati and P. Pfleiderer, "A theory of intraday patterns: Volume and price variability," *Review of Financial Studies*, vol. 1, pp. 3–40, 01 1988.
- [13] T. G. ANDERSEN, "Return volatility and trading volume: An information flow interpretation of stochastic volatility," *The Journal of Finance*, vol. 51, pp. 169–204, 03 1996.
- [14] L. Menkhoff, "The use of technical analysis by fund managers: International evidence," www.econstor.eu, 2010. [Online]. Available: <https://www.econstor.eu/handle/10419/38748>
- [15] M. S. Brown and M. Pelosi, "Moving averages trading method applied to cryptocurrencies," www.semanticscholar.org, 2019. [Online]. Available: <https://www.semanticscholar.org/paper/MOVING-AVERAGES-TRADING-METHOD-APPLIED-TO-Brown-Pelosi/81fad98931b58359a15f24b7bcada34ab/5973268>
- [16] F. Fang, C. Ventre, M. Basios, L. Kanthan, D. Martinez-Rego, F. Wu, and L. Li, "Cryptocurrency trading: a comprehensive survey," *Financial Innovation*, vol. 8, 02 2022.
- [17] Interdax, "Research: A variable moving average strategy for bitcoin outperforms hodling," Interdax Blog, 12 2019. [Online]. Available: <https://medium.com/interdax/research-a-variable-moving-average-strategy-for-bitcoin-outperforms-hodling-def78b27d8eb>
- [18] E. Kuo, B. Iles, and M. R. Cruz, "Ampleforth: A new synthetic commodity," *Ampleforth White Paper*, 2019.
- [19] E. Kuo, "The ampleforth + chainlink oracle integration is going live," Medium, 03 2020. [Online]. Available: <https://blog.ampleforth.org/the-ampleforth-chainlink-oracle-integration-is-going-live-16053ccdebd5>
- [20] L. Zhang and Y. Liu, "Optimal algorithmic monetary policy," *arXiv preprint arXiv:2104.07888*, 2021.
- [21] Y. Liu and L. Zhang, "Cryptocurrency valuation: An explainable ai approach," *arXiv preprint arXiv:2201.12893*, 2022.
- [22] L. Zhang and Z. Zhuang, "CryptoEnv: An automated trading environment for scientific research and applications in *Cryptoeconomics*," 2022. [Online]. Available: <https://github.com/sciecon/cryptoenv>