

Crypto Fraud Detection

Can-Elia Barth^{1*}, Florian Baumgartner^{1*} and Aaron Brülisauer^{1*}

¹FHNW University of Applied Sciences and Arts Northwestern
Switzerland.

*Corresponding author(s). E-mail(s): canelian.barth@students.fhnw.ch;
florian.baumgartner1@students.fhnw.ch;
aaron.brulisauer@students.fhnw.ch;

Abstract

Purpose: This project aims to detect cryptocurrency fraud by analyzing social media activity and financial data, focusing on the patterns preceding the collapse of scam coins. The primary objective is to identify features that differentiate scam coins from non-scam coins using machine learning techniques.

Methods: A dataset comprising social media posts from X/Twitter ($\sim 19k$ Posts) and Reddit ($\sim 450k$ Posts and Comments), combined with financial time series data, was created for both scam and non-scam coins. Features were extracted from embeddings of posts, price trends, and trading volume. These features were used to train machine learning models, including MultiRocket, to predict scam coins.

Results: The model showed high overall precision (0.78) in detecting scams but struggled with recall (0.60) with an accuracy of 0.60 and F1 Score of 0.53. Scam coins were more accurately identified when the scam was imminent, and the predictive performance decreased with longer time windows or unrelated historical data. Despite these limitations, the model demonstrated potential as an investor alert tool.

Conclusion: The results indicate that recent social media and financial activity contain critical features for identifying scams. Future work should focus on expanding the dataset and exploring additional deep learning methods to improve recall and generalization across a broader set of coins.

Keywords: Cryptocurrencies, Bitcoin, Ethereum, Fraud, Web Data Acquisition, Social Media, Reddit, X (Twitter), Time Series Analysis, MultiRocket, LSTM, Deep Learning

1 Introduction

The cryptocurrency market has evolved into a significant financial sector, but faces substantial challenges with fraudulent activities. Recent blockchain analysis shows that illicit cryptocurrency transactions summed up to \$24.2 billion in 2023 ([Chainalysis, 2024](#)). This broader trend of criminal activity is further evidenced by the Federal Trade Commission’s findings that cryptocurrency scams resulted in over \$1 billion in reported losses from January 2021 to March 2022, affecting more than 46,000 individuals, with investment schemes being the predominant form of fraud ([Federal Trade Commission, 2022](#)).

Unlike traditional financial markets, the cryptocurrency ecosystem operates in a largely unregulated environment, making it difficult to identify and prosecute fraudulent actors ([Mazorra, Adan, & Daza, 2022](#)). The inherent characteristics of cryptocurrencies—decentralization, anonymity, and ease of access—facilitate various types of fraud, including Ponzi schemes, pump-and-dump schemes, and rug pulls ([Scharfman, 2024](#)). These schemes exploit investor trust through aggressive promotion, where significant price increases often reduce investors’ willingness to scrutinize underlying activities for potential fraud ([Mazorra et al., 2022](#)).

This project does not categorise fraud types and does not focus on specific types of fraud. This project uses a broad approach: it tries to detect the price declines of fraudulent coins to near 0 ahead of time. The goal is to learn from well-known fraudulent coins by comparing them to legit coins. Specifically, it investigates how financial data and social media interactions can be used to predict the price fall and explores the potential of machine learning techniques for identifying scam coins. The research addresses the following questions:

1. How can financial data and social media interactions be utilized to detect fraudulent activities in the cryptocurrency market?
2. Which machine learning methods are most effective in identifying potential scam coins based on the collected data?

1.1 Related Work

Previous research has examined various aspects of cryptocurrency fraud. Studies have focused on transaction graph analysis and anomaly detection ([Li et al., 2019](#); [Patel, Pan, & Rajasegarar, 2020](#)), vulnerabilities in smart contracts ([Feist, Grieco, & Groce, 2019](#)), and automated scam detection mechanisms ([Mazorra et al., 2022](#)). However, these studies often focus on individual forms of fraud or rely solely on financial data. There is limited research integrating multiple data sources to comprehensively detect fraudulent activities ([Mazorra et al., 2022](#)). This project does not focus on rug pulls as [Mazorra et al. \(2022\)](#), but compares the price and social media data of well-known fraudulent coins and legit coins.

1.2 Preliminaries

1.2.1 Crypto Exchanges

Cryptocurrency exchanges are online platforms that facilitate the buying, selling, and trading of digital assets such as Bitcoin, Ethereum, and other cryptocurrencies. They serve as an essential gateway for individuals and institutions to participate in the crypto economy, offering functionalities similar to traditional stock exchanges but tailored to digital currencies (Ehrlich, 2023). These exchanges play a crucial role in establishing market liquidity and providing price discovery mechanisms.

There are two primary types of cryptocurrency exchanges: centralized exchanges (CEXs) and decentralized exchanges (DEXs). CEXs operate as intermediaries, managing trades and custody of funds on behalf of users. They often provide features like advanced trading tools, fiat-to-crypto transactions, and higher liquidity. However, their centralized nature makes them susceptible to hacking and fraud risks (Ehrlich, 2023).

In contrast, DEXs allow users to trade directly with one another using blockchain-based smart contracts. DEXs offer enhanced privacy and reduced reliance on intermediaries, but they may have limitations such as lower liquidity and a steeper learning curve for users (Ehrlich, 2023).

Despite their benefits, cryptocurrency exchanges have also been targets of fraud, including hacking, phishing, and exit scams. For example, the Bitfinex exchange suffered a breach in 2016, resulting in the theft of approximately 119,756 bitcoins (DOJ, 2024). Similarly, in 2022, the U.S. Department of Justice charged six individuals in four cases involving cryptocurrency-related fraud, with intended losses exceeding \$100 million (DOJ, 2022). Understanding the operational structure and security measures of these platforms is vital in developing robust detection mechanisms for crypto-related fraud.

1.2.2 Embedding

Word embeddings are numerical representations of text that map words or sentences into a high-dimensional vector space (Mikolov, Chen, Corrado, & Dean, 2013). These dense vector representations capture semantic relationships, where similar words have similar vector representations (Pennington, Socher, & Manning, 2014). The effectiveness of embeddings stems from the distributional hypothesis, which suggests that words appearing in similar contexts tend to have similar meanings (Harris, 1954). This property makes embeddings particularly useful for machine learning applications, as they transform textual data into a format that algorithms can process while preserving semantic relationships (Devlin, Chang, Lee, & Toutanova, 2019).

1.2.3 MultiRocket

MultiRocket, introduced by Tan, Dempster, Bergmeir, and Webb (2022), is a machine learning method for analyzing time series data. It extends the Rocket algorithm by Dempster, Schmidt, and Webb (2021) to handle multivariate time series (datasets with multiple features over time). By applying random convolutional kernels to the

data, MultiRocket generates features that simple models, like Random Forest, can use for classification. This method is fast and efficient because it transforms complex data into a format easier for traditional models to process. MultiRocket is particularly useful for detecting patterns in large datasets, such as identifying fraudulent activities in cryptocurrency markets.

2 Methods

To address the challenge of detecting fraudulent cryptocurrencies, a systematic approach was developed. This involved leveraging diverse data sources, including social media activity and financial metrics, and applying advanced machine learning techniques. The following section describes the steps taken to collect and preprocess the data, create meaningful features, and implement models for scam detection.

2.1 Data Collection

The general data collection approach involved gathering both social media and financial data to create a comprehensive dataset for each cryptocurrency. Social media data, including posts and comments, was scraped from X and Reddit using relevant keywords and hashtags related to cryptocurrencies. Financial data, such as price and trading volume, was retrieved in 4-hour intervals from [TradingView \(2025\)](#) to capture market activity.

2.1.1 Reddit Posts/Comments

Collecting posts and comments from Reddit proved to be a challenging task for this project. Initially, two different approaches were tested; however, the first approach had to be abandoned due to limitations on retrieving historical posts directly from the Reddit website.

In the first approach, an attempt was made to scrape posts from the official Reddit site using Selenium. However, the modern Reddit website quickly blocks remote access, leading to frequent connection issues. A possible workaround involved using the older version of Reddit ([old.reddit.com](#)), which is less restrictive. Nevertheless, this older site only allows approximately 250 posts to be loaded per query in a given subreddit, making it impossible to retrieve much older entries. The high volume of requests also triggered blocks on [old.reddit.com](#), which then necessitated the use of proxies.

Due to these constraints, a second approach was adopted. This alternative pipeline used a custom “Google-Scraper,” which searches Google for specific queries and aggregates corresponding Reddit links. For each collected link, the full content of the post and its comments were then retrieved separately through an application programming interface (API). This method enabled a more comprehensive retrieval of both posts and comments without being restricted by Reddit’s direct scraping limitations.

Figure 1 and Figure 2 illustrate the number of collected posts and comments per subreddit, respectively. It can be observed that the majority of the data includes “Bitcoin” as a term. Given Bitcoin’s status as the best-known cryptocurrency, a higher volume of discussion around it is to be expected.

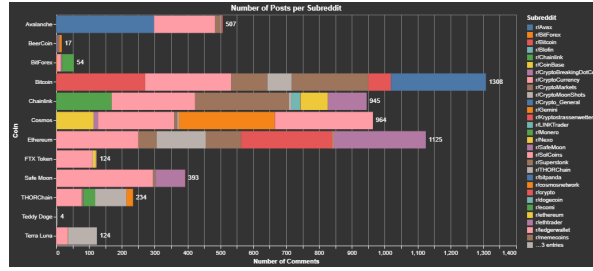


Fig. 1 Visualization of the Number of Posts per Subreddit and Term

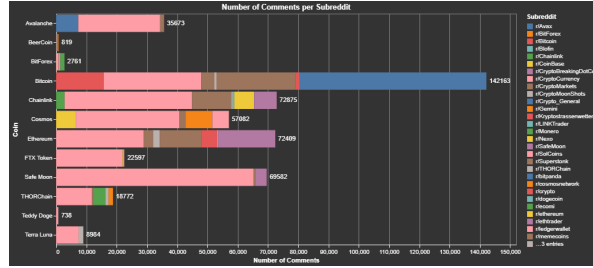


Fig. 2 Visualization of the Number of Comments per Subreddit and Term

In Figure 3, the number of Reddit posts per coin over time (after 2020) is presented. A notable trend is the increasing volume of posts as the timeline approaches the present day. This pattern can be attributed to the time-based nature of the “Google Scraper,” which inherently returns more recent content. Consequently, a higher number of posts in the most recent periods is to be expected.

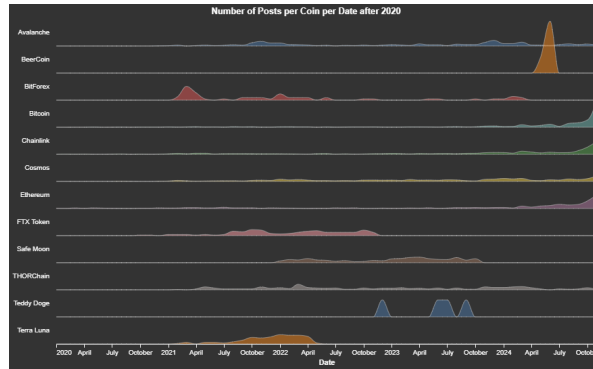


Fig. 3 Visualization of the Number of Posts per Coin after 2020

2.1.2 X (Twitter) Posts

The process began with data collection, where tweets related to cryptocurrencies were scraped using tools like Selenium to handle dynamically loaded content. Specific keywords, such as coin names, were used to filter relevant tweets. The scraped data included the tweet text, metadata (such as timestamps, likes, retweets, and comments), and user name.

In the following Figure are some examples of Coins of the collected posts from X (Figure 4) per Coin visualized.

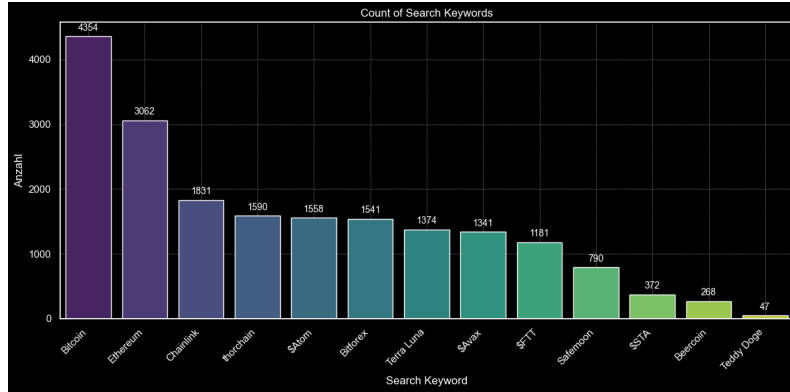


Fig. 4 Visualization of the Number of posts scraped from X per Coin

As seen in the figure above, the most collected data is same as in the reddit data from Bitcoin.

In Plot 5 there is an overview of the scraped Posts per Coin over time. In red are the start-date (date of our first price data point) and the end-date (date of scam if scam-coin or date of last price data point if non-scam coin).

The spikes, such as those at the beginning for Bitcoin, originated from a different implementation of the scraping algorithm that was used initially but later replaced with a more refined version. The data from the earlier implementation was not removed, as it still provides additional useful information and contributes to the overall dataset.

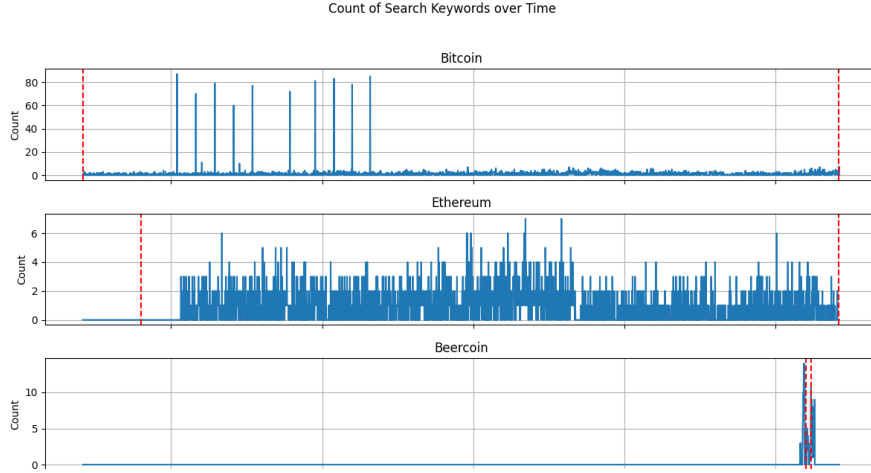


Fig. 5 Visualization of the posts scraped from X per Coin over time

The research process encountered several challenges during the data collection phase, particularly due to limitations and restrictions imposed by X (formerly Twitter) on scraping activity. These obstacles significantly impacted the speed and scale of data acquisition, ultimately limiting the size and diversity of the dataset.

One major issue was the need to mimic human behaviour during the scraping process to avoid detection by X's anti-bot mechanisms. To achieve this, numerous sleep intervals were incorporated. While this approach reduced the risk of being flagged as a bot, it significantly slowed down the scraping process, making it time-consuming to collect even a small volume of data.

Another challenge was the strict rate limiting enforced by X, which imposed a hard cap after approximately 200 scraped tweets. Once this limit was reached, a mandatory 10-minute pause had to be observed before resuming scraping, which also slowed down the scraping process.

Given the time-intensive nature of scraping under these constraints, the number of coins that could be included in the dataset was severely restricted. Leading to a dataset that, while rich in detail for the coins included (about 19'000 Posts over 13 Coins), was limited in its breadth.

2.1.3 Price Data

Price data and associated metrics, such as trading volume, were collected in 4-hour intervals from [TradingView \(2025\)](#) for all coins included in the dataset. This data was crucial for two purposes: first, to identify the date when scam coins experienced a significant price collapse for labelling purposes, and second, to extract additional features that could be incorporated into the machine learning models like opening, closing, high, and low prices, as well as trading volume.

2.2 Coin Labeling

This table shows which scam coins were used for the dataset.

Coin Name	Source	Reference	Scam Date
Terra Luna	SEC.gov	SEC (2023)	08.05.2022
FTX Token	Securities.io	Securities.io (2023)	07.11.2022
Beercoin	Krypto News	Krypto News (2024)	24.06.2024
BitForex	Crypto.news	Crypto.news (2024)	15.06.2024
SafeMoon	Reuters	Reuters (2023)	31.10.2023
Teddy Doge	CryptoSlate	CryptoSlate (2022)	21.07.2022
STA	Livemint	Livemint (2022)	17.07.2022

Table 1 Summary of Crypto Scams and Associated News with References

The labelling of coins as fraud or non-fraud was conducted based on two primary criteria: official news reports and market behaviour. Coins were labelled as fraud if credible news sources or announcements confirmed the project to be a scam or fraudulent operation. Additionally, this was cross-validated with market data, where coins that experienced a significant price collapse—falling to near-zero levels and failing to recover—were also classified as fraudulent. This dual approach ensured that the labelling process was both reliable and grounded in verifiable evidence.

2.3 Data Pipeline and Feature Creation

After downloading the posts, comments, and price data, several preprocessing steps were performed to prepare the data for modelling. These steps transformed the raw data into a structured, multivariate time series dataset for each coin, ready for use in machine learning models.

First, each post was embedded into a 512-dimensional vector using a pre-trained embedding model from Huggingface [Jina AI \(2023\)](#). These embeddings captured the semantic meaning of the posts. These embeddings were aligned with the price data by matching their timestamps to the corresponding 4-hour intervals. If multiple posts occurred within the same 4-hour tick, their embeddings were averaged to create a single representative vector for that interval. Conversely, if no posts were available for a given interval, a zero vector was inserted as a placeholder to maintain consistency in the time series. This process created a time series for each dimension of the 512-dimensional embedding vector, effectively treating each dimension as an independent time-series feature.

Next, additional time series data from the price and volume metrics was incorporated. Features such as opening price, closing price, high/low prices, and trading volume, all sampled at consistent intervals (e.g., 4 hours), were added to the dataset. Furthermore, more features like sentiment of social media posts was conducted on the textual data, and were included as additional time-series features. In [Figure 6](#) the data and model pipeline is visualized.

2.4 Scam Detection with Machine Learning

The preparation of multivariate time series data for training and evaluation followed a structured approach using fixed time windows. For each coin, multiple time windows were created based on random generated lengths and starting points. These windows were specifically designed to capture different temporal snapshots leading up to the scam date.

The time windows were defined in days and converted into corresponding data points (e.g., one day consisting of six data points for 4-hour intervals). Randomly generated time window lengths were capped at a maximum of 365 days, and starting points were fixed at specific offsets from 0 to 10 days before the scam date. For each combination of a time window and starting point, a padded time series was created for all features, ensuring consistency in shape across all data points.

To handle variations in the length of valid data available for a given coin and time window, left-padding was used for shorter series, aligning all time series to a maximum length equal to the longest time window. This padding ensured compatibility with models like MultiRocket, which require consistent input shapes.

The main model consisted of a pipeline that used MultiRocket from SKtime ([sktime \(2025\)](#)) as a feature extractor to combine the multivariate time series of a data point to one feature vector and classified the generated feature vector using a Random Forest classifier. The whole data and model pipeline is visualized in Figure 6.

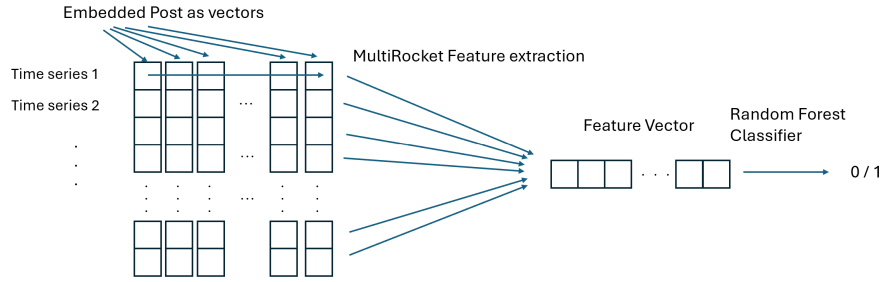


Fig. 6 Visualization of the Data and Model Pipeline

For the Test Data the same approach was used but only with an offset of maximum 5 Days before the scam. In the Train Set, it was set to 10 days to make the model more robust, but since the most important features were expected to appear closer to the scam date, the focus was on the last 5 days before the scam. The amount of input data was limited to avoid overwhelming the model with too many irrelevant features and ensure the focus remained on the most critical time periods leading up to the scam.

3 Results

3.1 Data Overview / EDA

Exploratory Data Analysis (EDA) focused primarily on the embeddings of the posts. These embeddings, representing the semantic content of the posts, were visualized in two dimensions using dimensionality reduction techniques like Principal Component Analysis (PCA) (Ringnér, 2008) and Uniform Manifold Approximation and Projection (UMAP) (McInnes, Healy, & Melville, 2018).

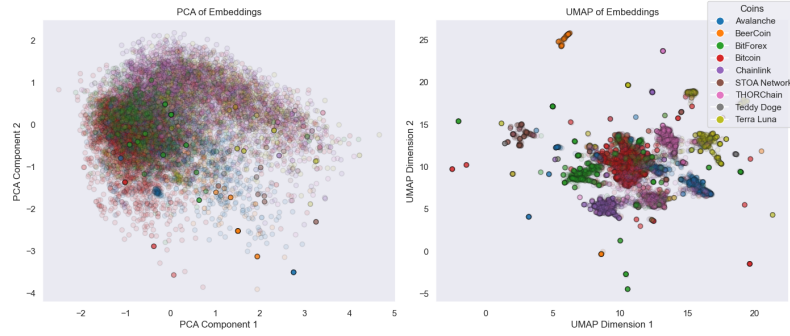


Fig. 7 Visualization of the posts from x embedded and plotted with PCA (left) and UMAP (right) with the Coins coloured

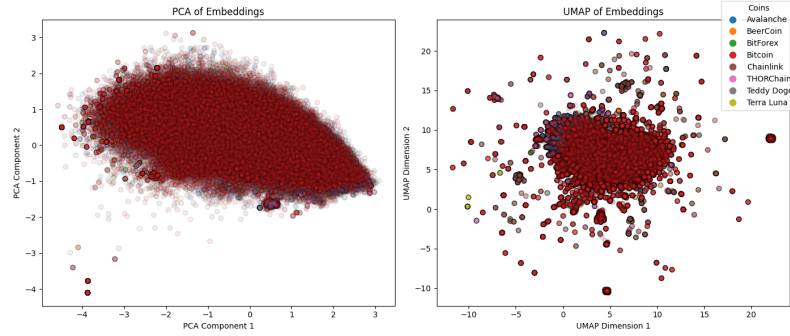


Fig. 8 Visualization of the posts from Reddit embedded and plotted with PCA (left) and UMAP (right) with the Coins coloured

It was observed that the embedding clusters of the individual coins exhibit distinct shapes and structures, especially in the X data. Further analyses were conducted, such as calculating metrics for each cluster to quantify their characteristics. These metrics provided additional insights into the differences between coins and their associated social media activity. All these analyses are documented in the provided code-repository.

3.2 ML-Performance on Test Set

Class	Precision	Recall	F1-Score	Support
0	0.56	1.00	0.72	100
1	1.00	0.21	0.35	100
Accuracy			0.60	200
Macro Avg	0.78	0.60	0.53	200
Weighted Avg	0.78	0.60	0.53	200

Table 2 Classification Report: Precision, Recall, F1-Score, and Support

For Class 0 (non-fraud coins), the model achieved a precision of 0.56, meaning that 56% of the predicted non-fraud instances were correct, and a recall of 1.00, indicating that all actual non-fraud coins and its data points were correctly identified. This resulted in an F1-score of 0.72.

For Class 1 (fraud coins), the precision was 1.00, suggesting that all instances predicted as fraud were indeed correct. However, the recall was much lower at 0.21, meaning that only 21% of the actual fraud coins were correctly identified. Consequently, the F1-score for this class was just 0.35, reflecting poor performance in detecting fraud due to the low recall.

Overall, the model achieved an accuracy of 60%, with a macro-average precision of 0.78, recall of 0.60, and F1-score of 0.53, while having a 50/50 class balance in the Test Set.

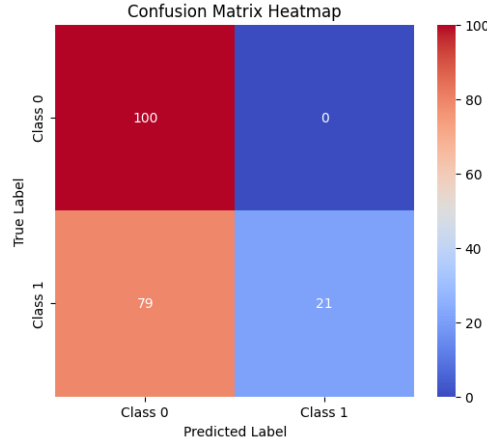


Fig. 9 Visualization of the Confusion Matrix of MultiRocket/Random Forest Pipeline Model on the Test Set

The model correctly classified 100 instances as Class 0, achieving perfect precision and recall for this class. However, it struggled significantly with identifying Class 1.

Out of 100 actual fraud instances, only 21 were correctly classified as fraud (true positives), while the remaining 79 were misclassified as non-fraud (false negatives). Notably, the model made no false positive errors, meaning it did not misclassify any non-fraud coins as fraud.

4 Discussion

The model’s performance is rather mixed. While it was able to perfectly predict non-scam coin windows, its performance on fraud coin windows was less reliable. However, when the model did predict a window as a scam, it was always correct. This indicates that the model has high precision for detecting scams but struggles with recall, as it often fails to identify all fraud coin windows, which means, that if the model predicts a fraud the probability of it being right is very high and should alarm investors. This trade-off highlights the model’s ability to confidently detect scams when enough evidence is present, but its limitations in consistently identifying all fraudulent activity.

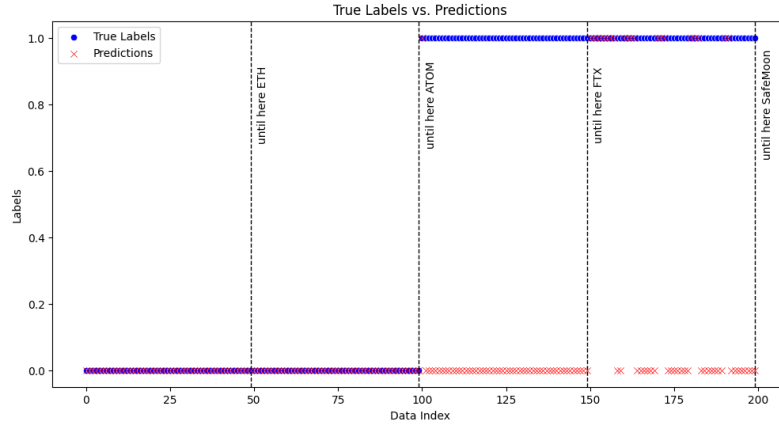


Fig. 10 Visualization of the real Labels and Predictions

The plot shows the true labels in blue, positioned at either 0 (non-scam) or 1 (scam), and the model’s predictions are marked with red crosses. The instances are grouped by coin, starting with Ethereum (ETH), followed by Cosmos (ATOM), FTX, and finally Safemoon. It is evident that the model predicts non-scam coins (ETH and ATOM) very well, with almost perfect alignment between the true labels and predictions.

For FTX, however, the model performs poorly, with many incorrect predictions and a lack of consistent alignment with the true labels. In contrast, for Safemoon, the model demonstrates a noticeable pattern in its predictions, where it is able to identify scam-related instances more effectively. This pattern was further analyzed and broken down to investigate how the individual instances were structured and why the model performed better in this specific case.

This could be due to the fact that Safemoon is a "typical" scam coin, heavily promoted on social media using bots, which makes the fraudulent activity more apparent in the embeddings of the posts. These patterns are likely easier for the model to detect.

In contrast, FTX does not fit the profile of a typical scam coin. The scam was primarily tied to the exchange behind the coin, rather than the coin itself. Additionally, FTX does not exhibit the same rug-pull behavior typically associated with scams, and it was not aggressively inflated through social media bot activity. As a result, the key features that indicate a scam in the embeddings are less prominent, making it more challenging for the model to accurately predict fraud in this case.

Despite these challenges, the model was able to detect FTX as a scam in one instance when using the shortest window without cutting off any of the newest days. This suggests that there are indeed some subtle features in the data that point toward fraudulent activity, even for a less typical scam coin like FTX. These features, while sparse, seem to provide just enough information for the model to make a correct prediction in certain scenarios.

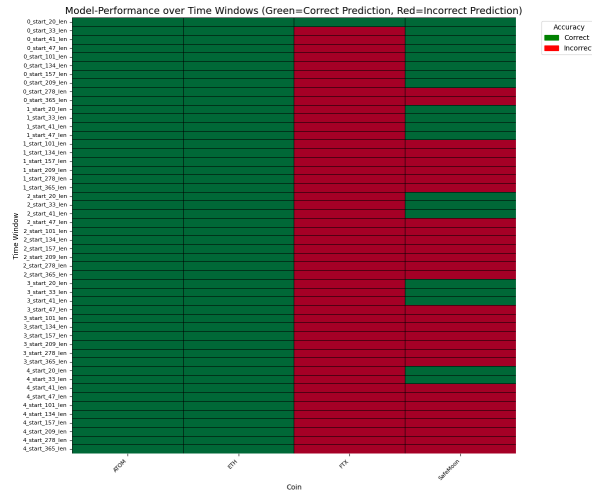


Fig. 11 Visualization of the Coin Instance Windows and the Modell Predictions

Figure 11 shows different coin windows on the left and the corresponding coins along the bottom. The colours indicate whether the model's predictions were correct (green) or incorrect (red). Each window is labeled with two key numbers: the first, labeled "start," represents the number of days before the scam date that were excluded from the model; the second, labeled "len", represents the length of the window in days.

For example, a label like 1.start_101.len refers to a window starting 1 day before the scam date and extending 101 days backward, using those features for the model.

The model appears to predict more instances correctly for Safemoon when more recent days are included in the window, i.e., when the x_start value is smaller. Conversely, as the window length increases, the model's performance deteriorates. This

suggests that the data from the days immediately preceding the scam contains crucial features indicating that the coin is a scam.

However, these features do not seem strong enough to dominate when the window becomes longer. In longer windows, the critical pre-scam features may be overshadowed by other features from much earlier in time, when the coin’s characteristics still resembled those of a legitimate project. Additionally, cutting off data from a few days before the scam further diminishes the importance of these key features, causing the model to make incorrect predictions more frequently as the window size increases. This highlights the importance of recent data in identifying scams effectively.

In summary, our model performs better at predicting a scam when the scam is imminent. The closer the scam date, the more accurate the model’s predictions become. Additionally, smaller windows yield better results, as it appears that scams are primarily detectable shortly before they occur, rather than over extended periods.

To utilize our model as a tool for investor protection, one could collect data from the last 20–30 days for a given coin and maybe a second time with even a smaller day range to not sparse out the important features, input this data into the model, and evaluate its prediction. If the model identifies a scam, this should trigger a high alert, prompting further investigation and caution. However, if the model does not detect a scam, this should not be taken as a guarantee of safety. This is particularly true for atypical scam coins, such as FTX, where the model may struggle to make accurate predictions due to unique characteristics that deviate from typical scam patterns.

4.1 Limitations

This project was limited in several ways. One significant limitation was the small number of coins included in the dataset due to the extensive effort required for scraping the data. This restriction reduces the statistical significance of the model and its results, as the limited dataset may not capture the full diversity of scam and non-scam coins.

Additionally, the model’s performance was weakened by the extreme noise present in the data. This issue is evident from instances where inputs containing the same recent days but additional days from further away from the scam date led to different and often incorrect predictions. This suggests that irrelevant or less important features from earlier periods diluted the impact of critical patterns in the days leading up to a scam, further challenging the model’s ability to generalize effectively.

4.2 Outlook

It would be crucial for the continuation of this project to significantly expand the dataset by including a larger number of coins. This would enable the development of more robust models and yield more statistically significant results. Testing this improvement could involve evaluating model performance on an expanded dataset to determine whether the increased diversity and size lead to better generalization and reliability in detecting scams.

4.2.1 LSTM

The LSTM model could not be finalized due to time constraints and the focus on the MultiRocket approach. To complete it, a more suitable implementation of the evaluation metrics would have been necessary to ensure proper performance assessment and comparison with other models.

4.2.2 Other Deep Learning Approaches

The ConvTran model (Foumani, Tan, Webb, and Salehi (2023)) was experimented with and successfully brought to training. However, it did not perform well, as the code was too encapsulated to be effectively adapted to our specific task. Additionally, the dataset was too small to adequately train such a complex model, further limiting its effectiveness. To further develop this approach, more time would need to be invested in adapting the implementation to better align with the specific task. Additionally, a significantly larger number of coins would need to be included in the dataset to provide sufficient data for training such a complex model effectively.

Acknowledgements

We, the authors, would like to extend our special thanks to Gabriel Torres Gamez for his collaboration with us on gathering and conducting exploratory data analysis (EDA) on social media data.

We also wish to express our gratitude to Moritz Kirschmann and Adrian Brändli for their invaluable guidance throughout this Challenge-X. Moritz especially supported us in developing machine learning strategies and guiding the process from raw data to a functional ML model. Adrian provided critical assistance in defining fraud within the cryptocurrency space and addressing the societal embedding of a data science research question within the broader societal context.

Additionally, we would like to thank Roman Studer and Cédric Huwyler for their insightful advice regarding our machine learning approach. Their expertise was instrumental in refining the strategies used in this project. We are also grateful to Fernando Benites for his guidance on web scraping, which played a key role in gathering the necessary data for our analysis.

Declarations

- AI tools
Parts of the text in this paper have been corrected for grammar and improved for readability with the assistance of ChatGPT (OpenAI, 2025).
- Code availability
All code is available on github.com/CryptoFraudDetection.
- Author contribution
All authors contributed equally to all aspects of this work.

References

- Chainalysis, t. (2024, January). *2024 Crypto Crime Trends from Chainalysis*. Retrieved 2025-01-17, from <https://www.chainalysis.com/blog/2024-crypto-crime-report-introduction/>
- Crypto.news (2024). *Crypto exchange bitforex left customers without \$55.6m in alleged exit scam*. Retrieved from <https://crypto.news/crypto-exchange-bitforex-left-customers-without-55-6m-in-alleged-exit-scam/> (Accessed: 2025-01-18)
- CryptoSlate (2022). *Teddy doge developers pull out \$4.5 million in alleged 'soft rug pull'*. Retrieved from <https://cryptoslate.com/teddy-doge-developers-pull-out-4-5-million-in-alleged-soft-rug-pull/> (Accessed: 2025-01-18)
- Dempster, A., Schmidt, D.F., Webb, G.I. (2021, August). MiniRocket: A Very Fast (Almost) Deterministic Transform for Time Series Classification. *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining* (pp. 248–257). Virtual Event Singapore: ACM. Retrieved 2025-01-17, from <https://dl.acm.org/doi/10.1145/3447548.3467231>
- Devlin, J., Chang, M.-W., Lee, K., Toutanova, K. (2019). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, ,
- DOJ (2022, June). *Office of Public Affairs | Justice Department Announces Enforcement Action Charging Six Individuals with Cryptocurrency Fraud Offenses in Cases Involving Over \$100 Million in Intended Losses | United States Department of Justice*. Retrieved 2025-01-17, from <https://www.justice.gov/opa/pr/justice-department-announces-enforcement-action-charging-six-individuals-cryptocurrency-fraud>
- DOJ (2024, November). *Office of Public Affairs | Bitfinex Hacker Sentenced in Money Laundering Conspiracy Involving Billions in Stolen Cryptocurrency | United States Department of Justice*. Retrieved 2025-01-17, from <https://www.justice.gov/opa/pr/bitfinex-hacker-sentenced-money-laundering-conspiracy-involving-billions-stolen>
- Ehrlich, S. (2023, May). *Crypto Exchanges: What Investors Need To Know*. Retrieved 2025-01-17, from <https://www.forbes.com/sites/digital-assets/article/crypto-exchanges-what-investors-need-to-know/> (Section: Forbes Digital Assets)
- Federal Trade Commission (2022, June). *Reports show scammers cashing in on crypto craze*. Retrieved from <https://www.ftc.gov/news-events/data-visualizations/data-spotlight/2022/06/reports-show-scammers-cashing-crypto-craze> (Accessed: 2025-01-17)

- Feist, J., Grieco, G., Groce, A. (2019). Slither: a static analysis framework for smart contracts. *2019 IEEE/ACM 2nd International Workshop on Emerging Trends in Software Engineering for Blockchain (Wetseb)*.
- Foumani, N.M., Tan, C.W., Webb, G.I., Salehi, M. (2023, September). Improving position encoding of transformers for multivariate time series classification. *Data Mining and Knowledge Discovery*, 38(1), , <https://doi.org/10.1007/s10618-023-00948-2> Retrieved from <http://dx.doi.org/10.1007/s10618-023-00948-2>
- Harris, Z.S. (1954). Distributional structure. *Word*, 10(2-3), 146–162,
- Jina AI (2023). *Jina embeddings v2 small english model*. Retrieved from <https://huggingface.co/jinaai/jina-embeddings-v2-small-en> (Accessed: 2025-01-18)
- Krypto News (2024). *Insiderhandel bei beercoin: Verdächtige kursmanipulationen und warnungen vor meme-coin scams*. Retrieved from <https://news-krypto.de/krypto-news/bitcoin/insiderhandel-bei-beercoin-verdaechtige-kursmanipulationen-und-warnungen-vor-meme-coin-scams/> (Accessed: 2025-01-18)
- Li, Y., Islambekov, U., Akcora, C., Smirnova, E., Gel, Y.R., Kantarcioglu, M. (2019, December). *Dissecting Ethereum Blockchain Analytics: What We Learn from Topology and Geometry of Ethereum Graph*. arXiv. Retrieved 2025-01-17, from <http://arxiv.org/abs/1912.10105> (arXiv:1912.10105 [cs])
- Livemint (2022). *Odisha police arrests india head of mega cryptocurrency scam worth rs 1,000 cr*. Retrieved from <https://www.livemint.com/news/india/rs-1-000-crore-cryptocurrency-scam-odisha-police-arrests-indias-head-of-the-fraudsters-networkbiharjharkhanddelhi-11691464062393.html> (Accessed: 2025-01-18)
- Mazorra, B., Adan, V., Daza, V. (2022, January). *Do not rug on me: Zero-dimensional Scam Detection*. arXiv. Retrieved 2024-09-24, from <http://arxiv.org/abs/2201.07220> (arXiv:2201.07220 [cs, q-fin])
- McInnes, L., Healy, J., Melville, J. (2018). Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, ,
- Mikolov, T., Chen, K., Corrado, G., Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, ,
- OpenAI (2025, January). *ChatGPT | OpenAI*. Retrieved 2025-01-17, from <https://openai.com/chatgpt/overview/>

- Patel, V., Pan, L., Rajasegarar, S. (2020). Graph deep learning based anomaly detection in ethereum blockchain network. *International conference on network and system security*.
- Pennington, J., Socher, R., Manning, C.D. (2014). Glove: Global vectors for word representation. *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 1532–1543,
- Reuters (2023). *Safemoon executives charged in us with fraud related to crypto token*. Retrieved from <https://www.reuters.com/legal/safemoon-executives-charged-us-with-fraud-related-crypto-token-2023-11-01> (Accessed: 2025-01-18)
- Ringnér, M. (2008). What is principal component analysis? *Nature biotechnology*, 26(3), 303–304,
- Scharfman, J. (2024). *The Cryptocurrency and Digital Asset Fraud Casebook, Volume II: DeFi, NFTs, DAOs, Meme Coins, and Other Digital Asset Hacks*. Cham: Springer Nature Switzerland. Retrieved 2024-11-07, from <https://link.springer.com/10.1007/978-3-031-60836-0>
- SEC (2023). *Sec charges terraform and ceo do kwon with defrauding investors in crypto schemes*. Retrieved from <https://www.sec.gov/newsroom/press-releases/2023-32> (Accessed: 2025-01-18)
- Securities.io (2023). *Sec and cftc charge co-conspirators in ftx collapse; ftt labeled "crypto security token"*. Retrieved from <https://www.securities.io/sec-and-cftc-charge-co-conspirators-in-ftx-collapse-fft-labeled-crypto-security-token/> (Accessed: 2025-01-18)
- sktime (2025). *Minirocketmultivariate - sktime v0.20.0*. Retrieved from https://www.sktime.net/en/v0.20.0/api_reference/auto_generated/sktime.transformations.panel.rocket.MiniRocketMultivariate.html (Accessed: 2025-01-18)
- Tan, C.W., Dempster, A., Bergmeir, C., Webb, G.I. (2022, February). *MultiRocket: Multiple pooling operators and transformations for fast and effective time series classification*. arXiv. Retrieved 2025-01-17, from <http://arxiv.org/abs/2102.00457> (arXiv:2102.00457 [cs])
- TradingView (2025). *Tradingview: Advanced financial charts & tools*. Retrieved from <https://www.tradingview.com> (Accessed: 2025-01-18)