

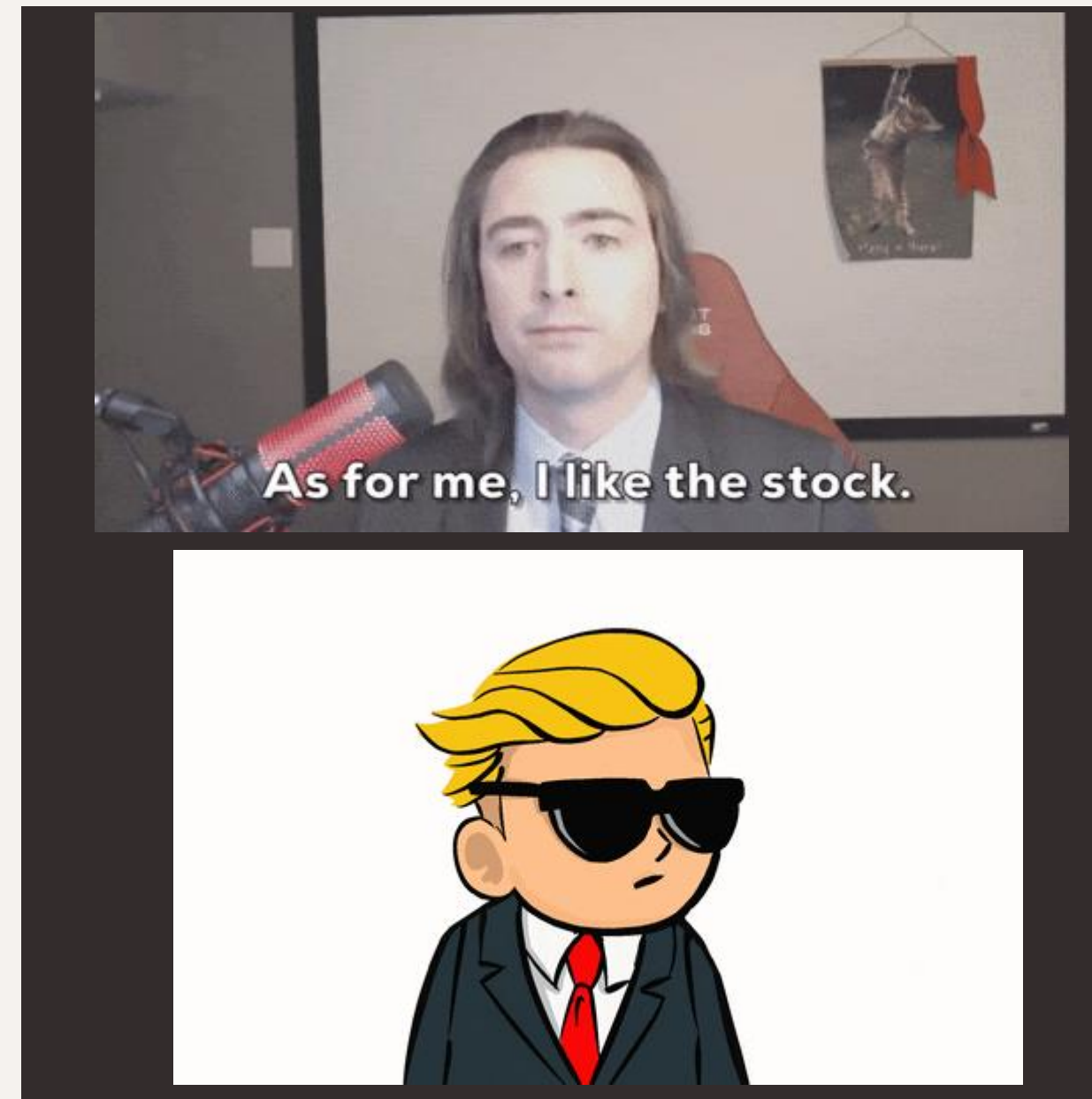
Riding the Waves of Crypto and Stocks: NLP-Powered Voyage through Reddit's Investment Chatter

Develop a binary NLP classifier to distinguish between crypto and stock posts found on the r/wallstreetbets and r/CryptoMoonShots subreddits. This will help users identify investment opportunities in both markets and potentially inform trading decisions based on which assets are popular/rising. #ShortSqueeze #TotheMoon



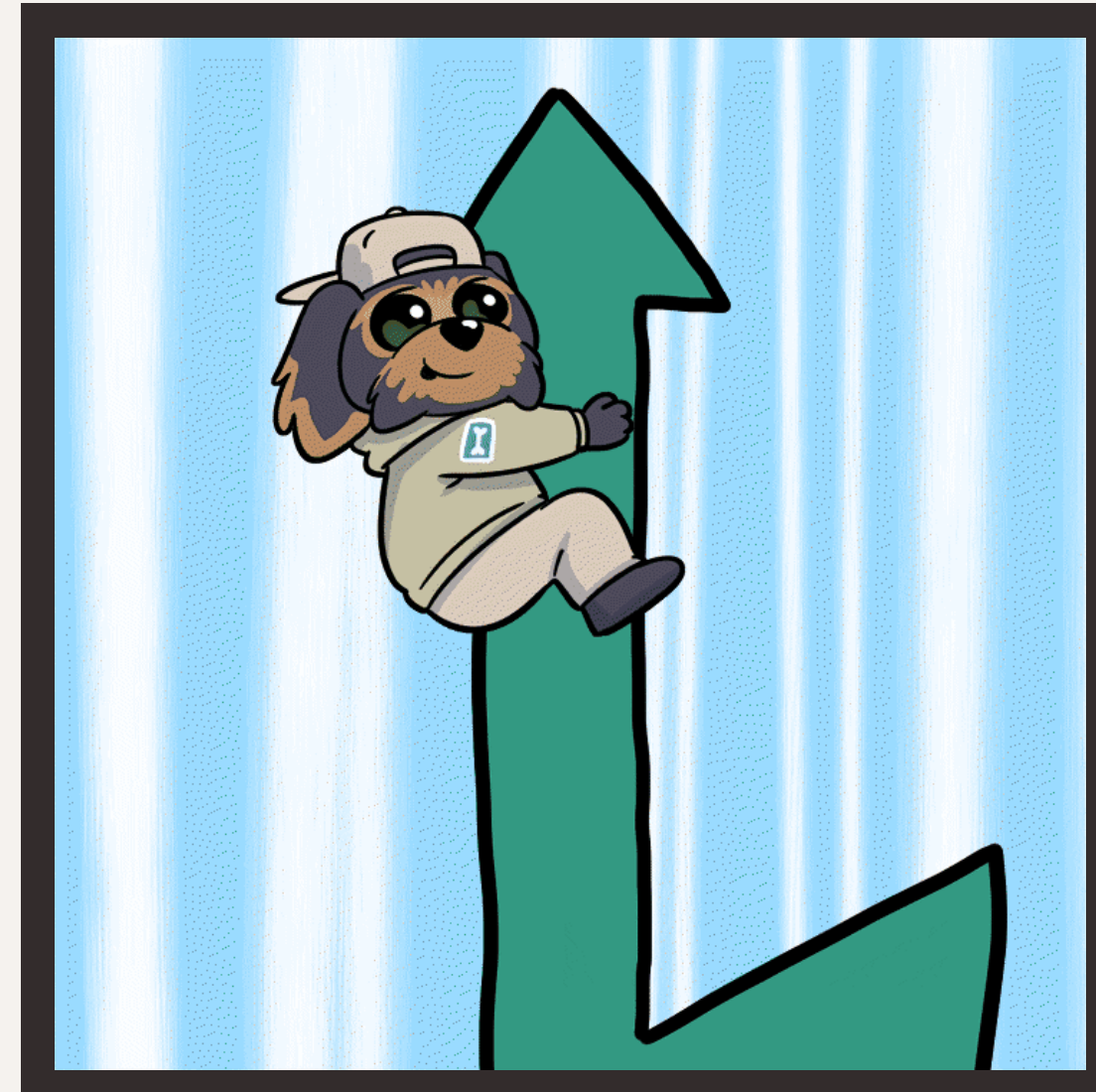
r/wallstreetbets

- **Meme Stonks & Retail Rebellion:** WallStreetBets rose to fame during the pandemic with the GameStop saga. Users coordinated massive stock purchases, popularizing the term "meme stocks" and challenging Wall Street institutions. This "casino-style" approach to the market has become a point of legal concern.
- **RoaringKitty (Keith Gill):** popularized the GameStop short squeeze in 2021, raising questions about retail investor influence and sparking a trial that could define legal boundaries for online financial communication.
- **Financial Lingo Goes Viral:** WSB's crass humor and unique financial slang ("diamond hands," "to the moon") became mainstream conversation in financial news.



r/CryptoMoonShots

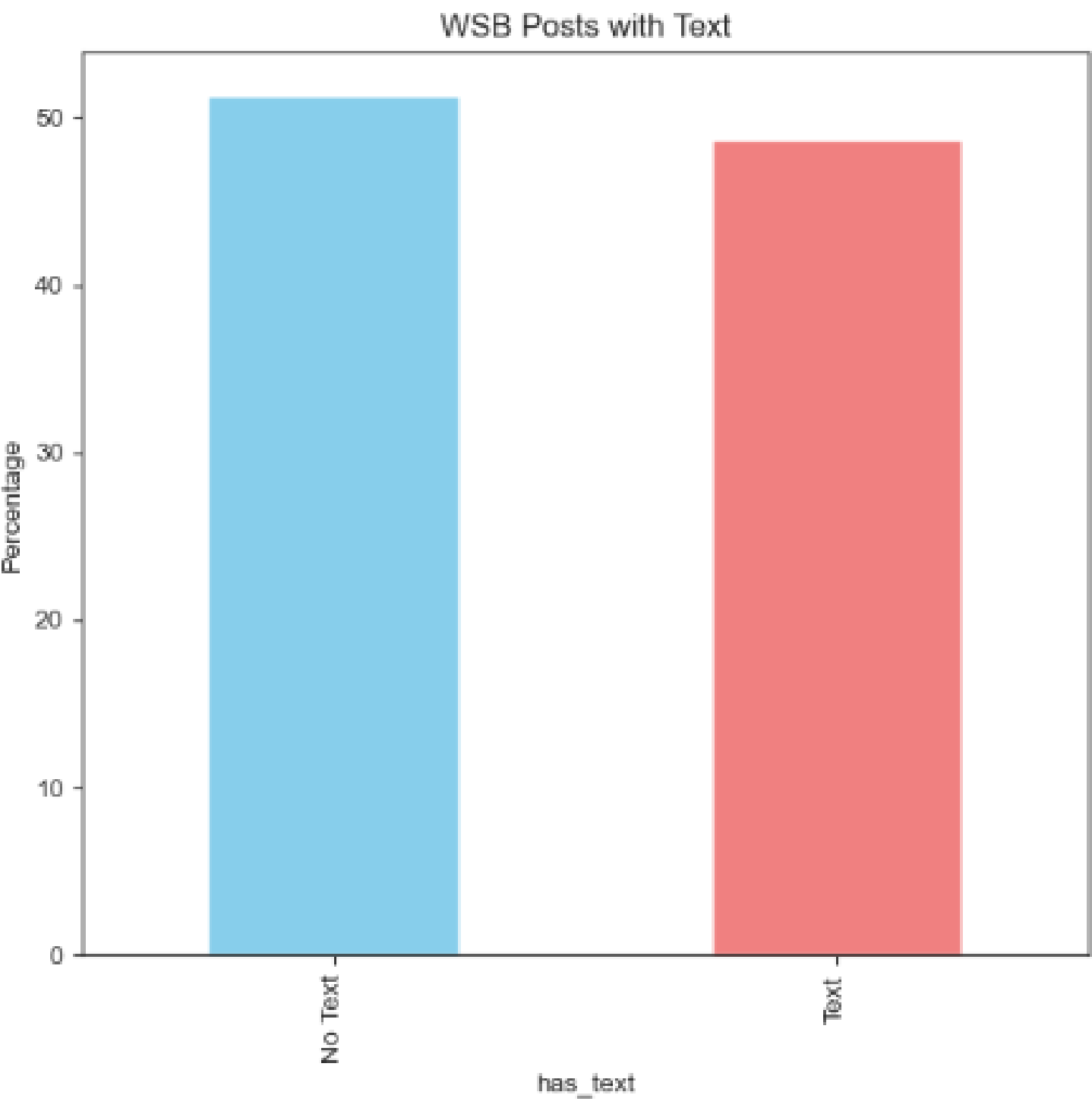
- **Community Shilling & High Risk:** The subreddit thrives on user recommendations and promotions (shilling) of unproven cryptocurrencies. This can lead to high-risk investments and frequent scams.
- **"To the Moon" & Memes:** Similar to WallStreetBets, CryptoMoonShots uses memes and internet slang ("to the moon," rocket emojis) to express excitement about crypto gains.
- **Focus on Price & Short-Term Gains:** Discussions often center around short-term price movements and potential profits, rather than the underlying technology or long-term viability of crypto projects.



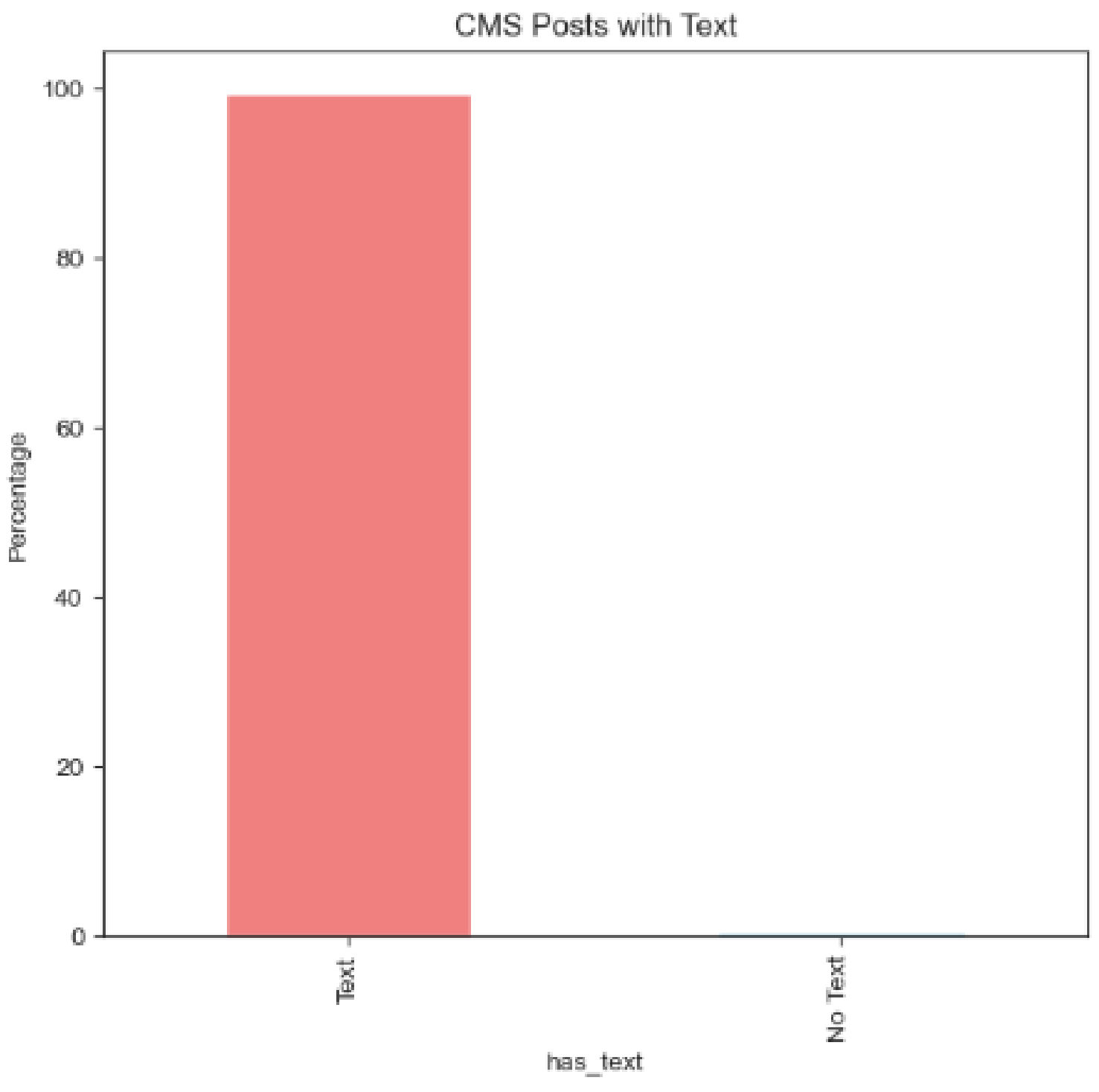
Objectives

- 1. Gather and prepare** data from the wallstreetbets and CryptoMoonShots subreddits using PRAW (Python Reddit API Wrapper). Extract from top, new and controversial sections of each subreddit.
- 2. Preprocess** the text data by cleaning, tokenizing, and vectorizing the posts for NLP analysis.
- 3. Train and compare** classification models – multinomial naïve bayes, random forest trees and logistic regression - to predict whether a post is related to cryptocurrency or stocks.
- 4. Evaluate** the performance of the models using appropriate metrics such **as accuracy, precision, recall, and F1 score**.

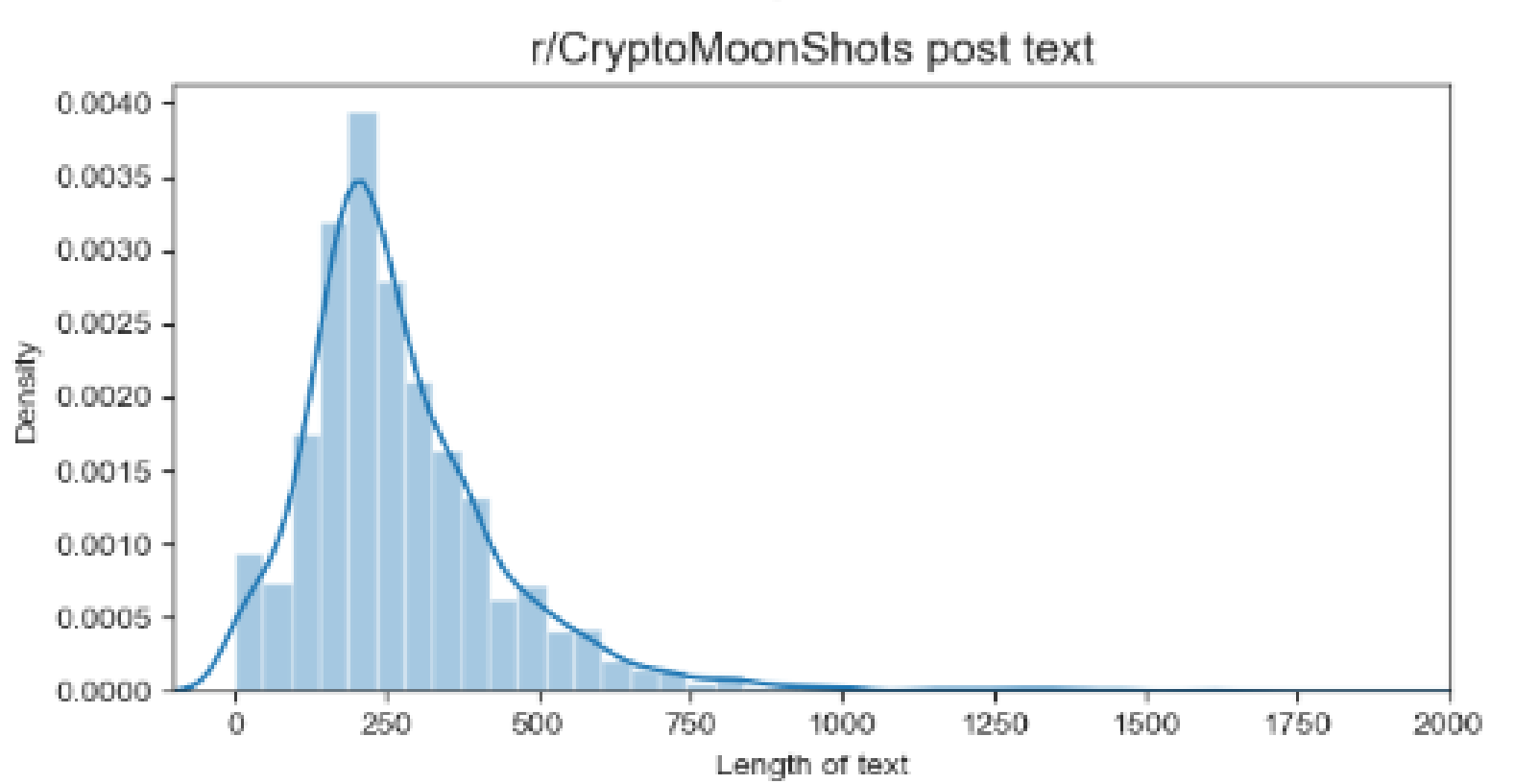
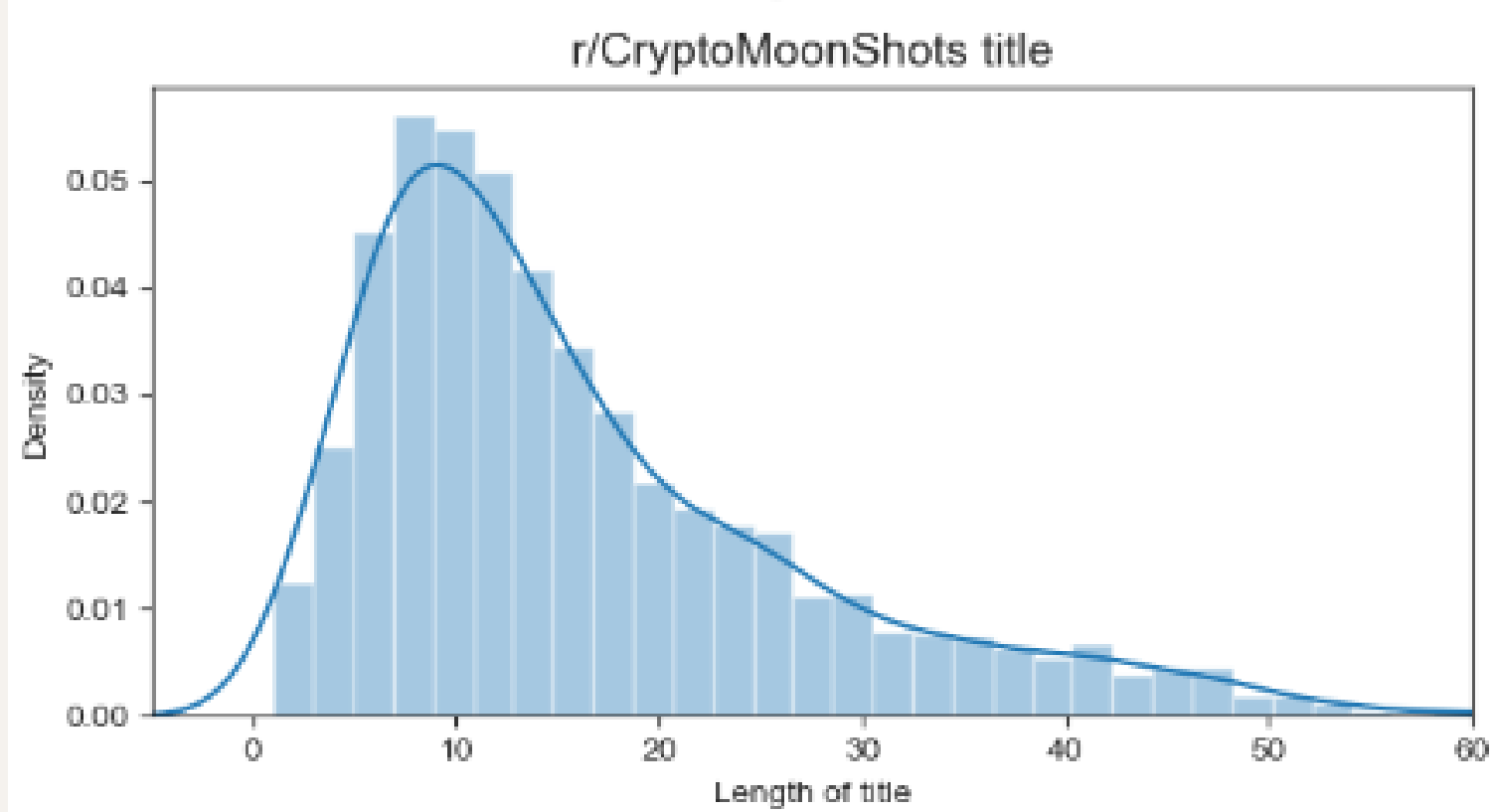
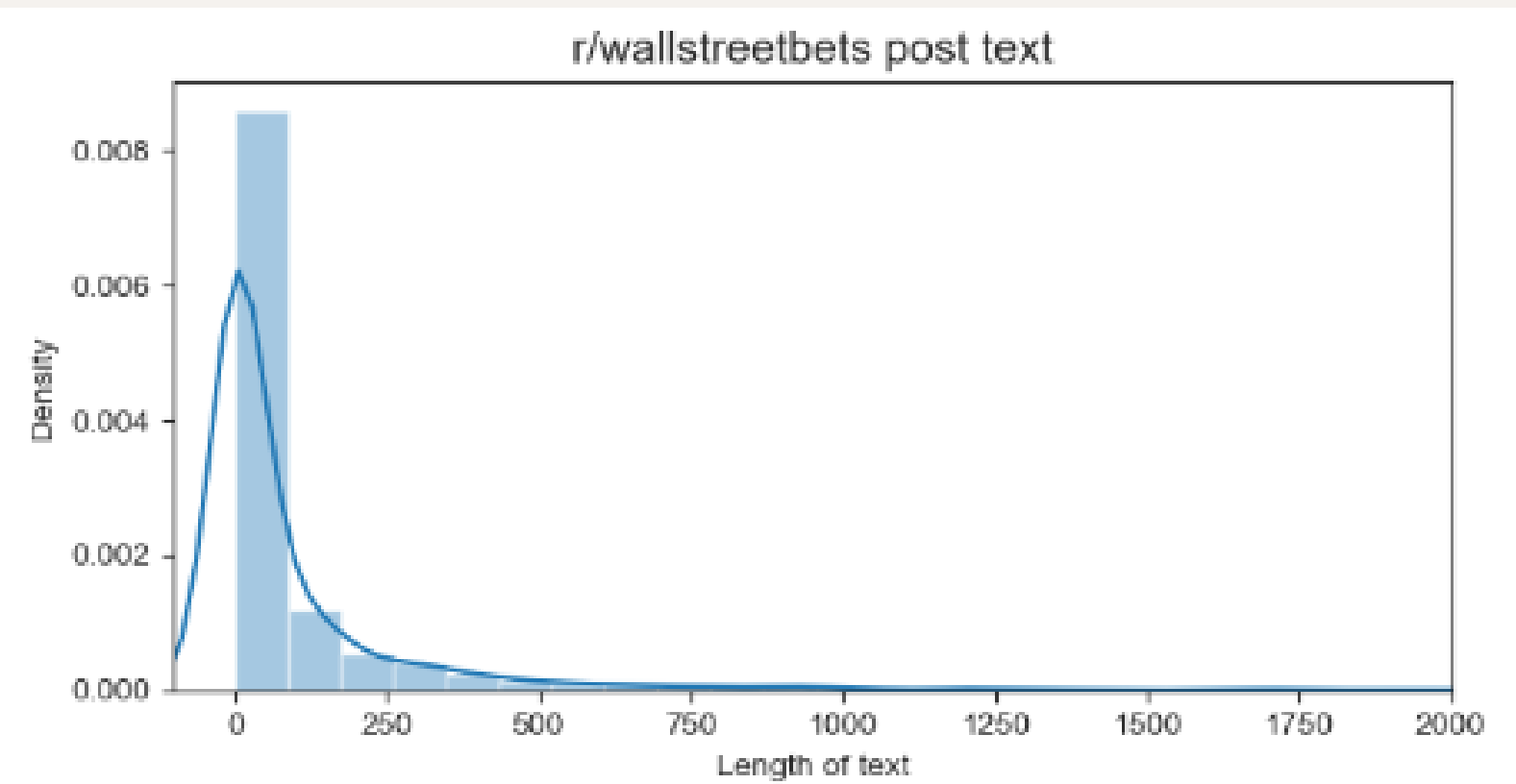
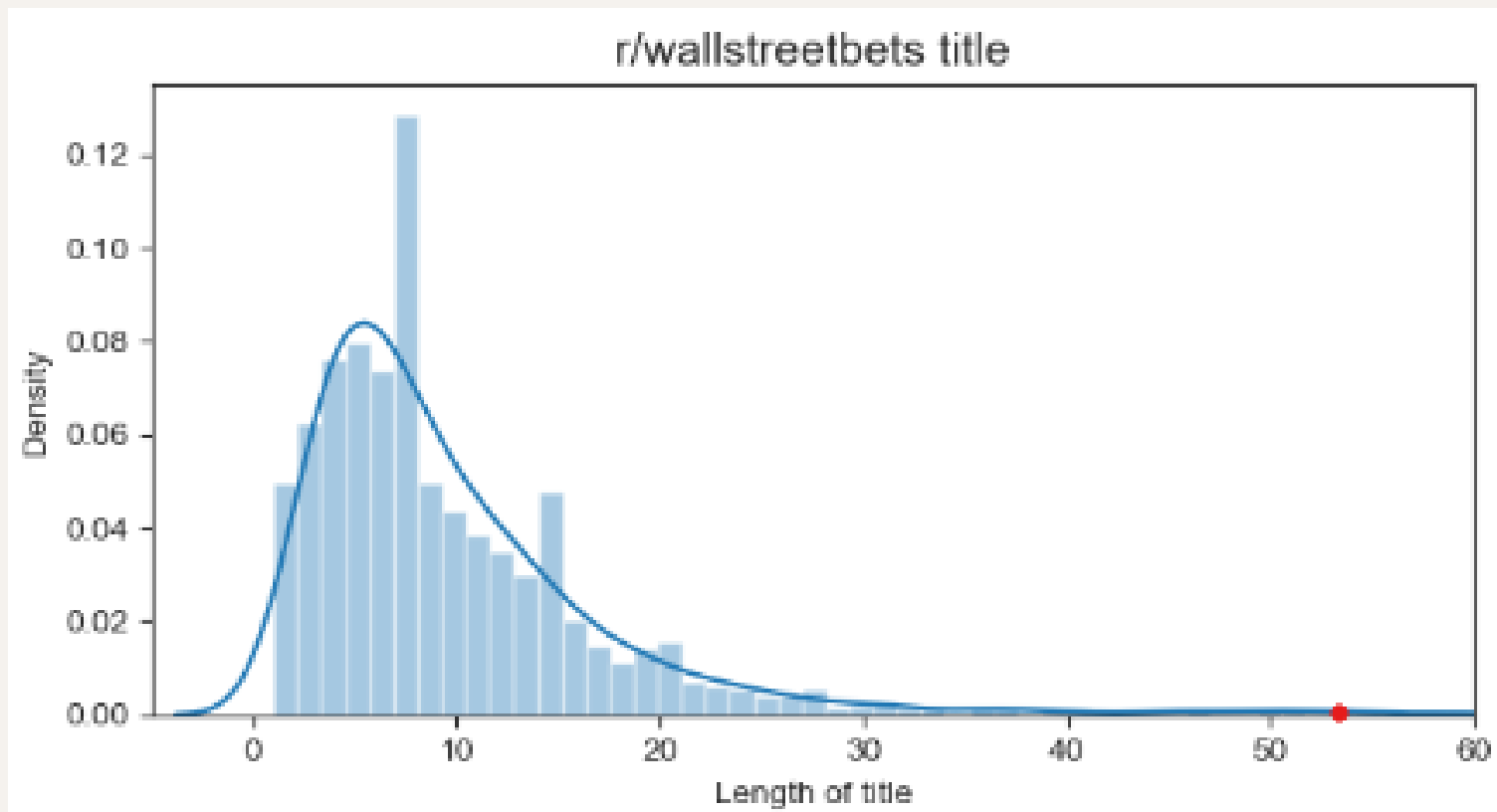


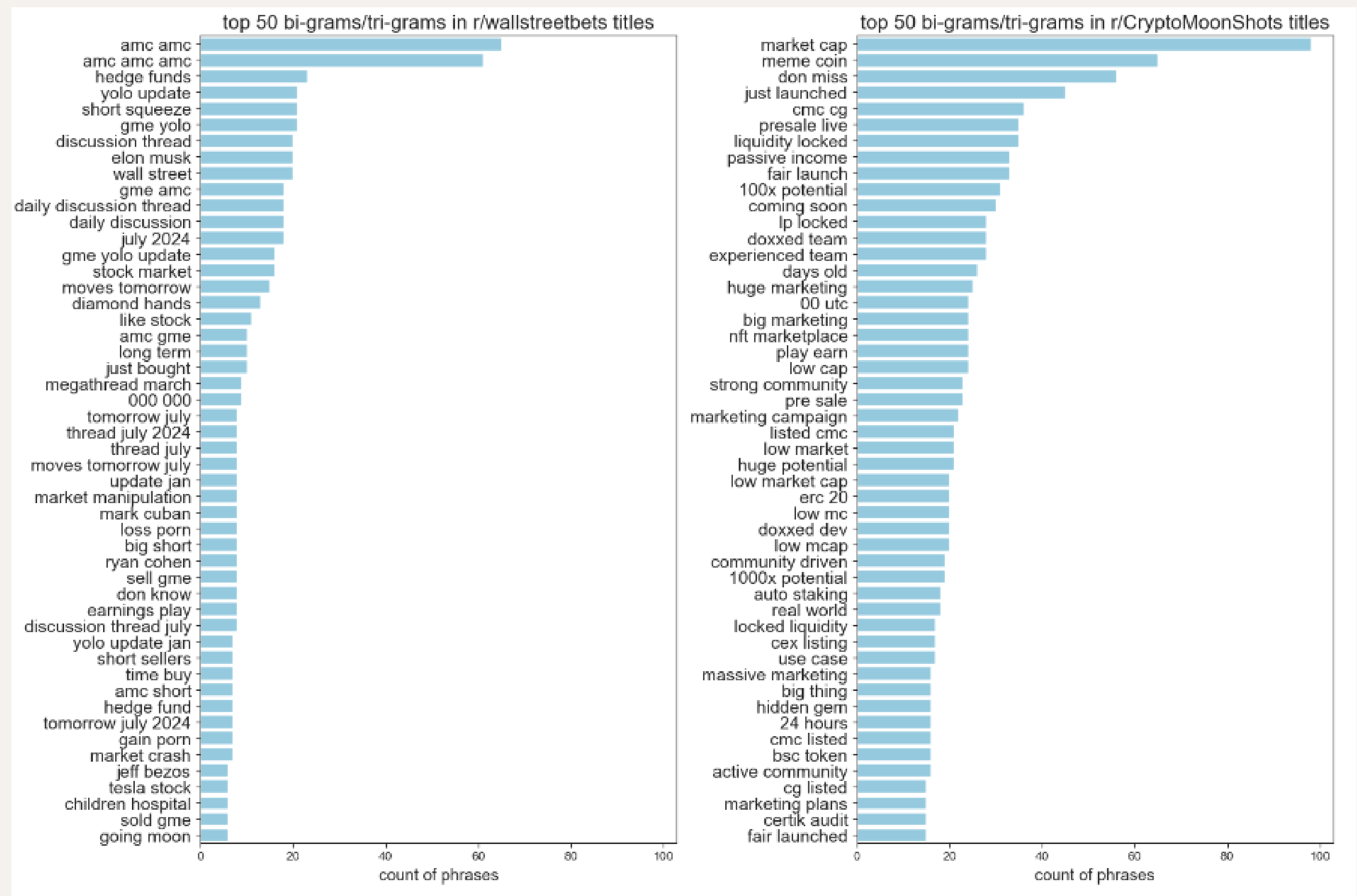


```
has_text
False    51.292951
True     48.707049
```



```
has_text
True     99.375459
False    0.624541
```



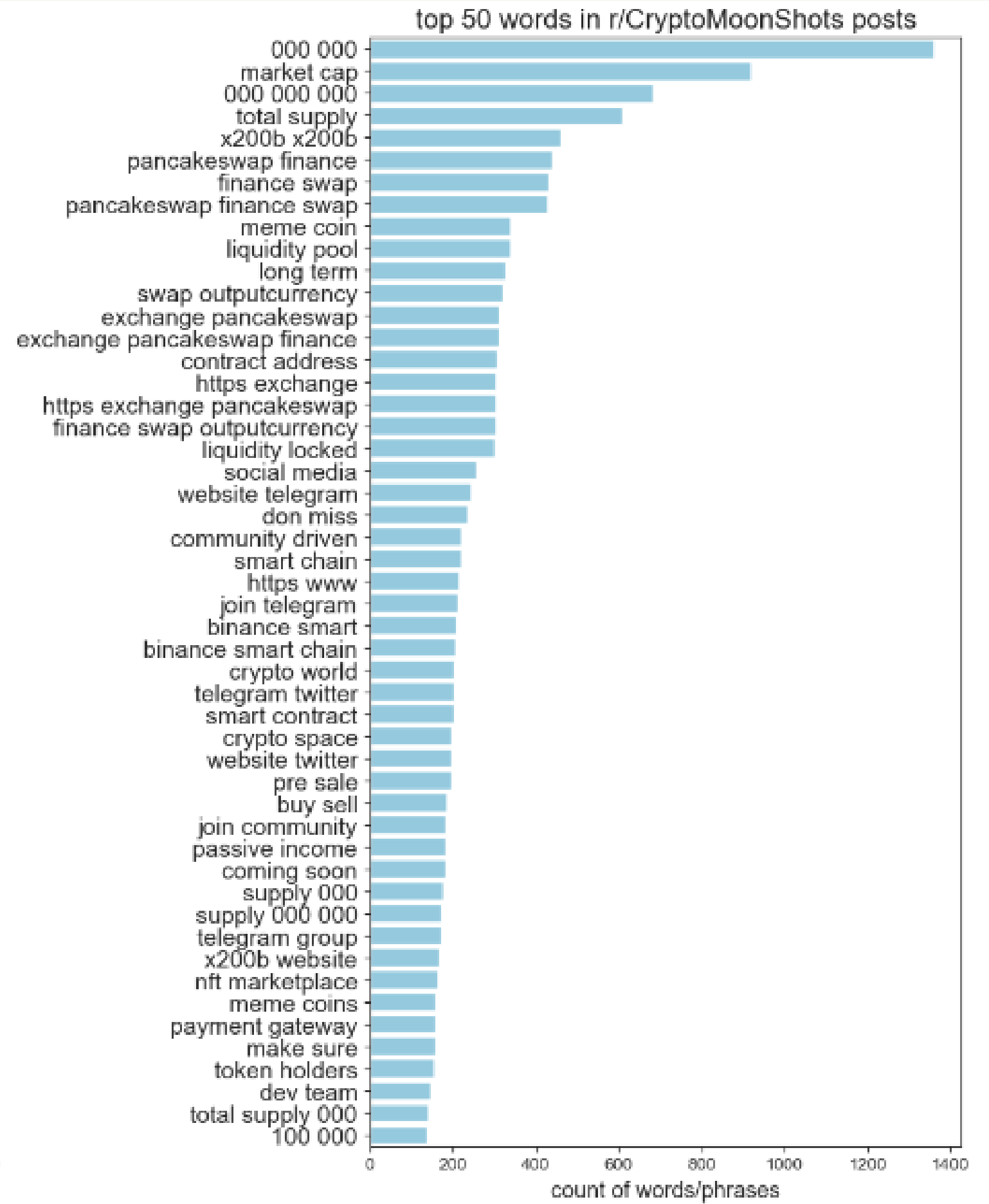
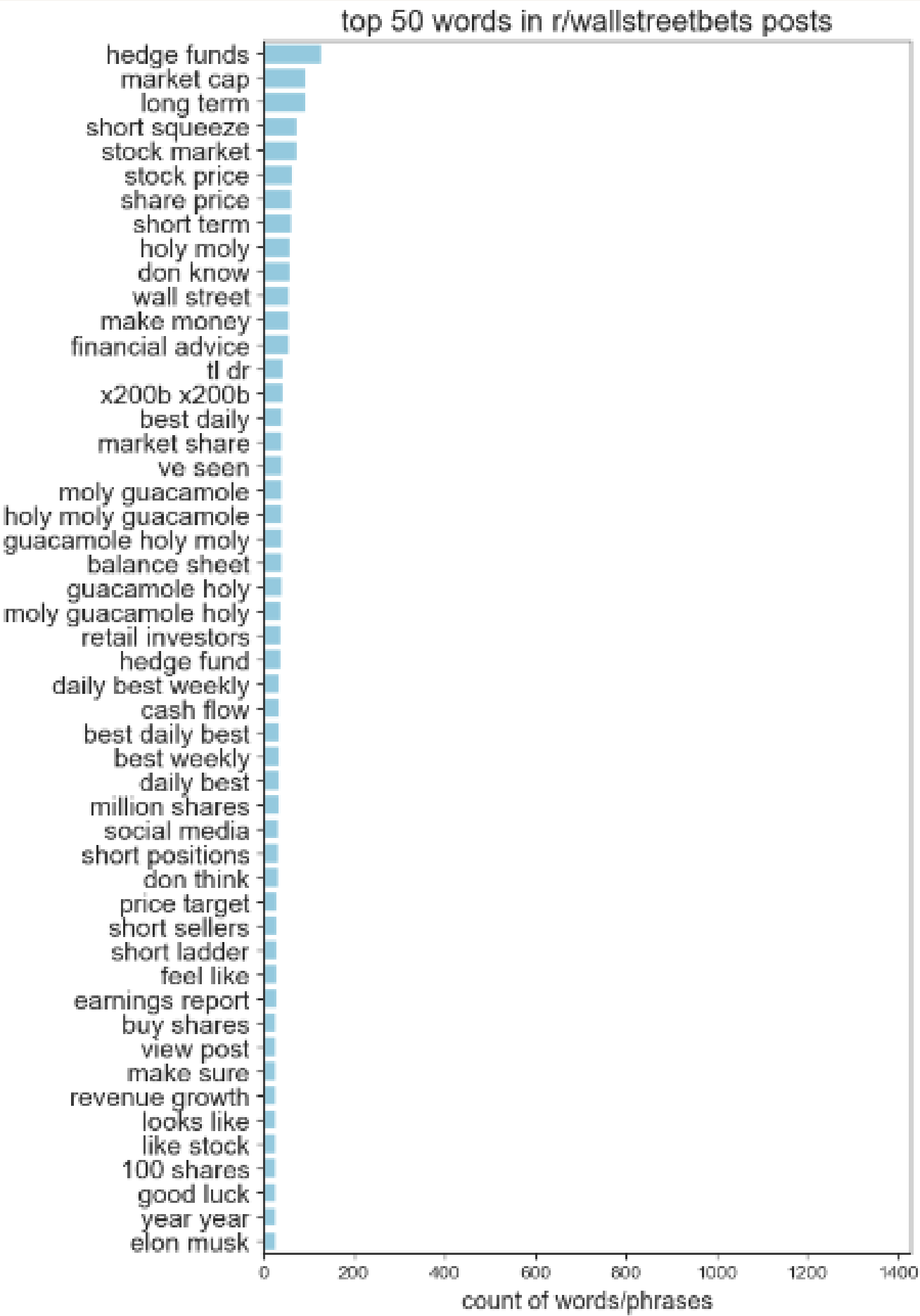


r/wallstreetbets

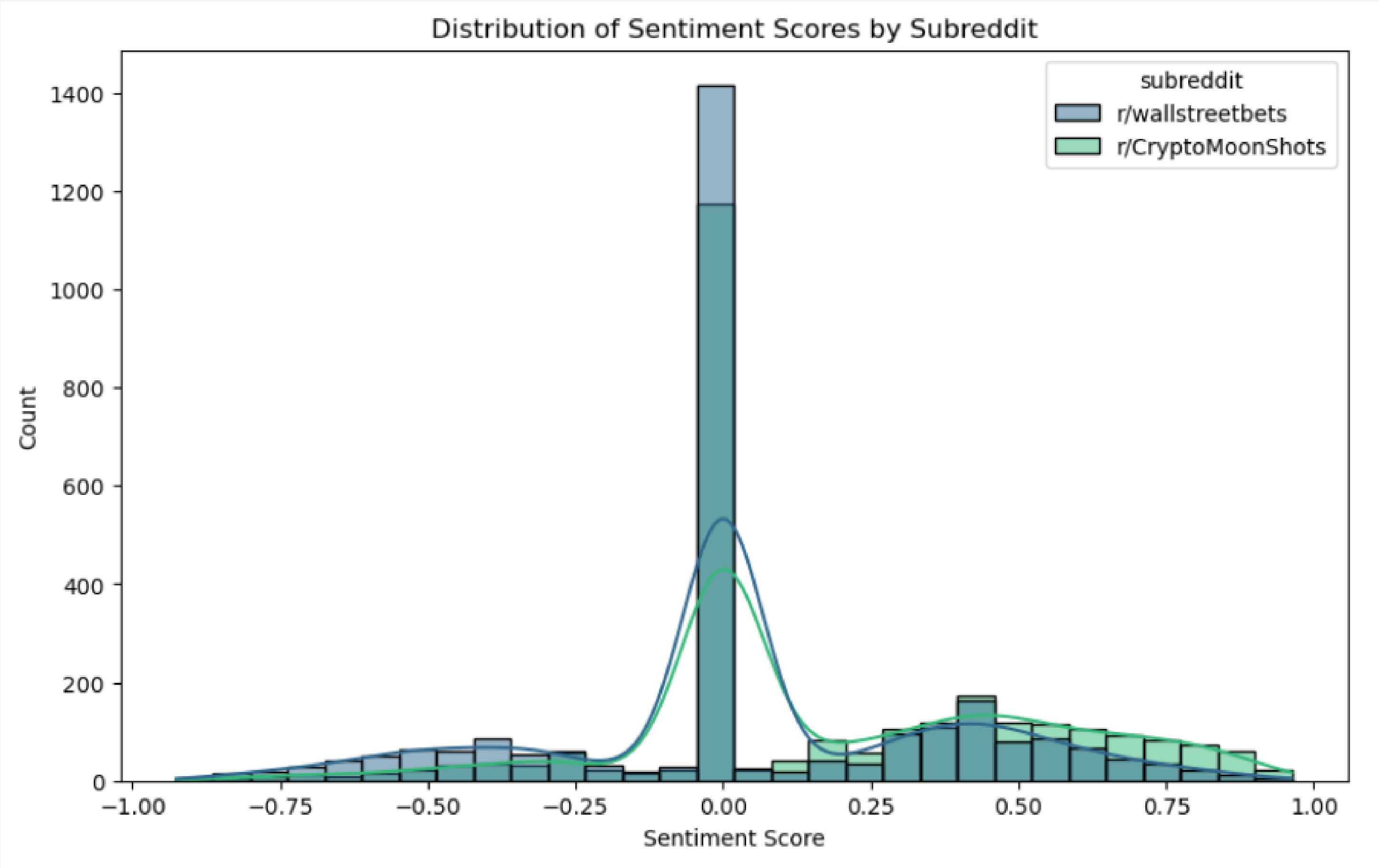
discussion thread
amc
elon musk
hedge funds
amc
short squeeze
gme
wall street
yolo update
diamond hands
don know
market manipulation
yolo update jan
amc short
children hospital
tomorrow july 2024
tomorrow july
earnings play
going moon
tesla stock
stock market
ryan cohen
thread july 2024
big short
just bought
daily discussion thread
update jan
market crash
update jan
gme yolo update
sold gme
jeff bezos
moves tomorrow
long term
sell gme
megathread march
mark cuban
time buy
000 000
discussion thread july
00 utc
nft marketplace
100x potential
experienced team
doxxed team
cmc cg
huge potential
low cap
fair launch
pre sale
1000x potential
lp locked
locked liquidity
pre sale
don miss
cmc cg
100x potential
days old
nft marketplace
00 utc
nft marketplace
days old
just launched
market cap
launched
just
market
presale live
meme coin
liquidity locked
cex listing
fair launched
marketing campaign
low mc
marketing plans
cg listed
low mcap
listed cmc
play earn
big marketing
income
active community
passive marketing
massive marketing

r/CryptoMoonShots

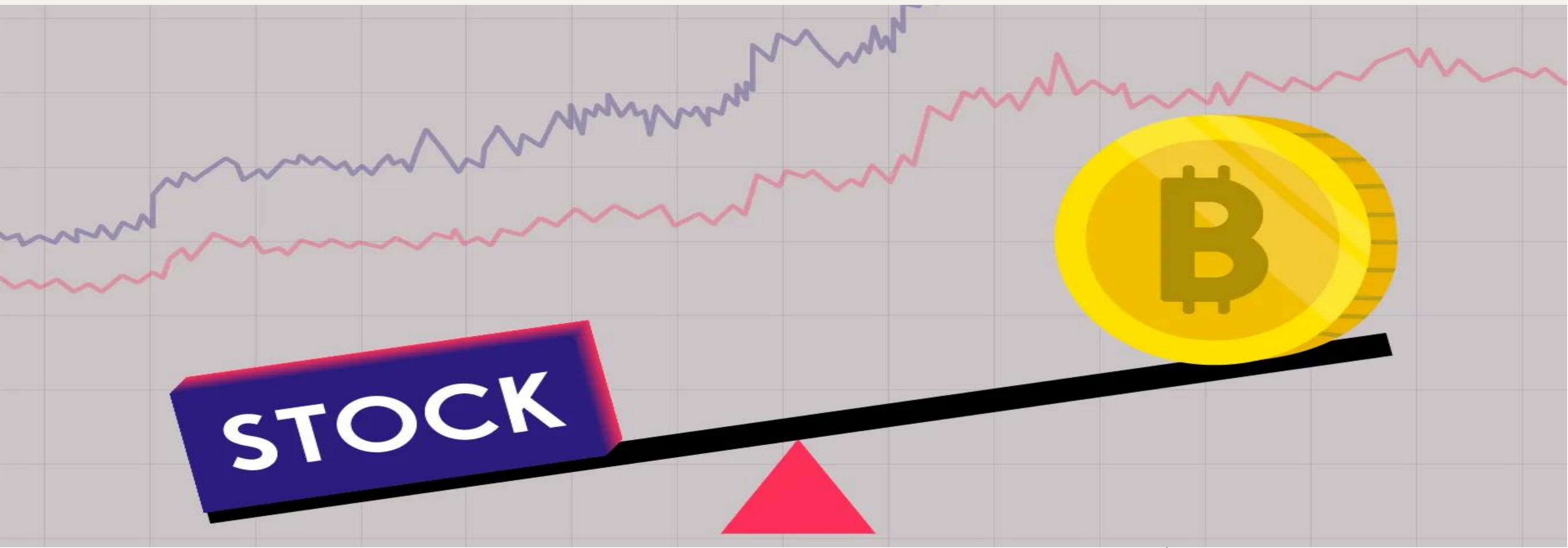
liquidity locked
meme coin
presale live
market cap
just launched
don miss
cmc cg
100x potential
days old
nft marketplace
00 utc
nft marketplace
days old
just launched
market cap
launched
just
market
presale live
meme coin
liquidity locked
cex listing
fair launched
marketing campaign
low mc
marketing plans
cg listed
low mcap
listed cmc
play earn
big marketing
income
active community
passive marketing
massive marketing



[illegible][illegible][illegible][illegible]



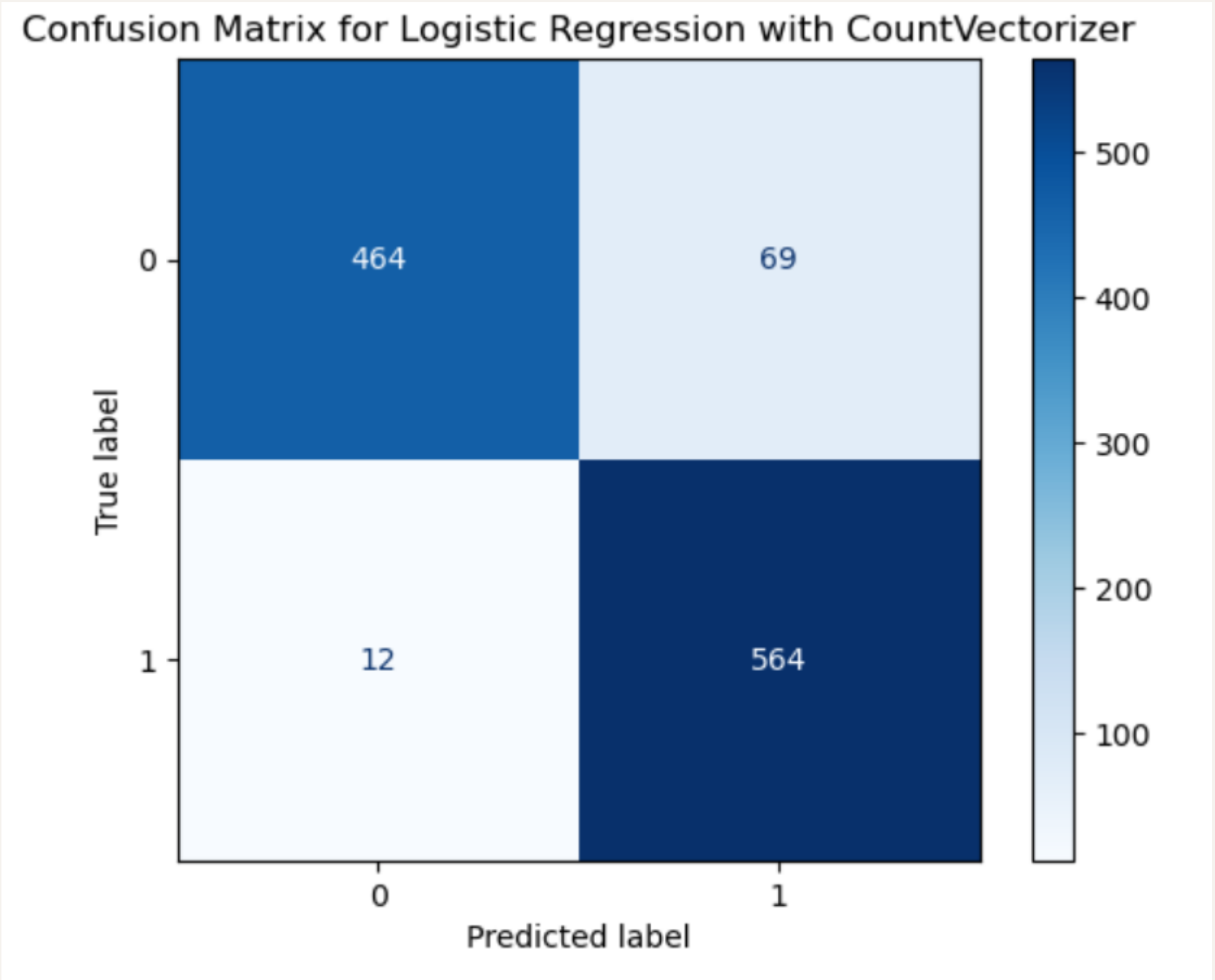
- After merging all dataframes into one combined csv file, mapping 0 to r/CryptoMoonShots and 1 to r/wallstreetbets, as well as tokenizing the title text into a processed feature for sentiment analysis, we can establish our baseline accuracy!
- **Baseline Accuracy – 0.519 or 51.9%**



Model 1: Logistic Regression with CountVectorizer

Training Accuracy	Testing Accuracy
0.928	0.927

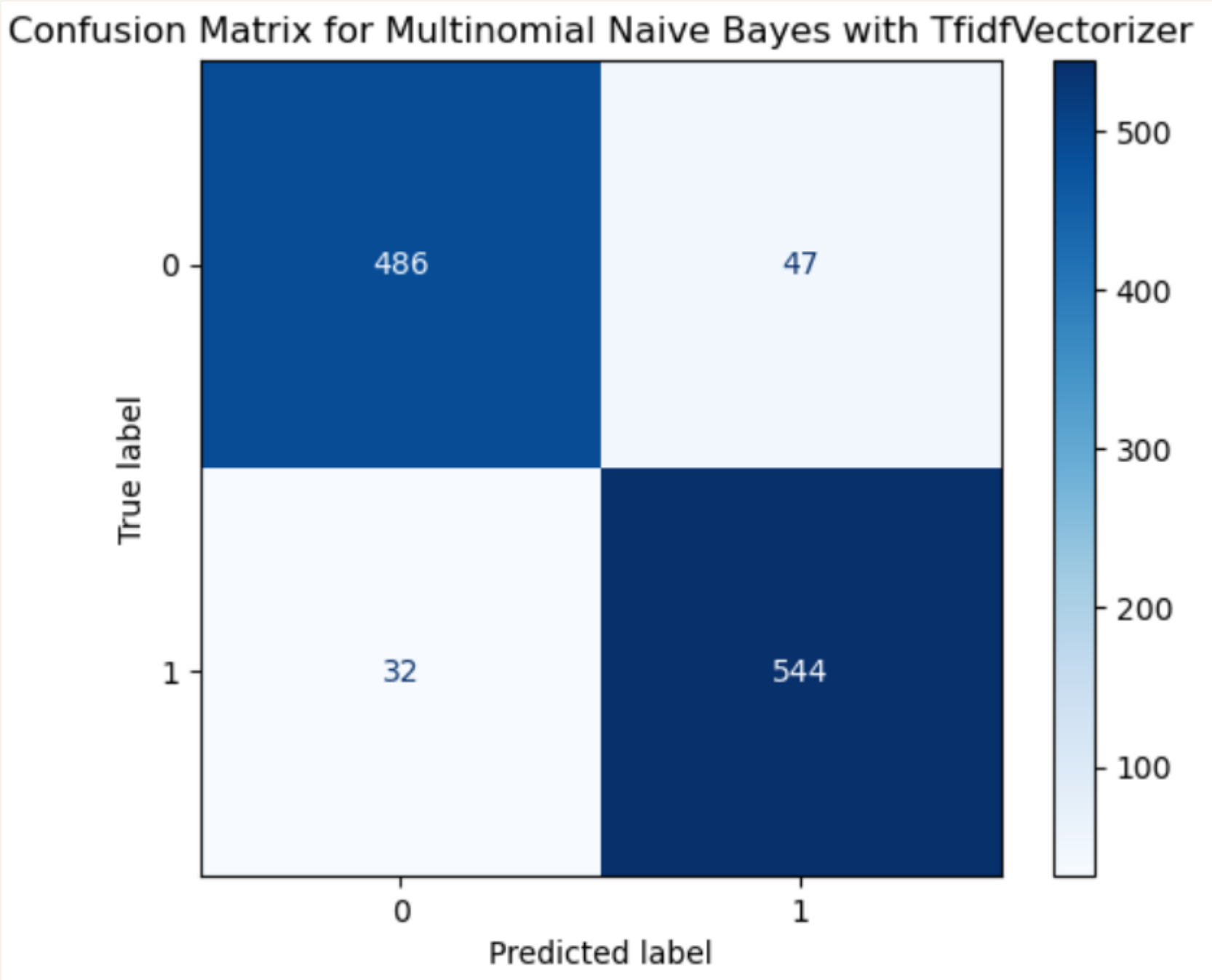
Classification Report for Logistic Regression:				
	precision	recall	f1-score	support
0	0.97	0.87	0.92	533
1	0.89	0.98	0.93	576
accuracy			0.93	1109
macro avg	0.93	0.92	0.93	1109
weighted avg	0.93	0.93	0.93	1109



Model 2: Multinomial Naïve Bayes with TfidfVectorizer

Training Accuracy	Testing Accuracy
0.928	0.929

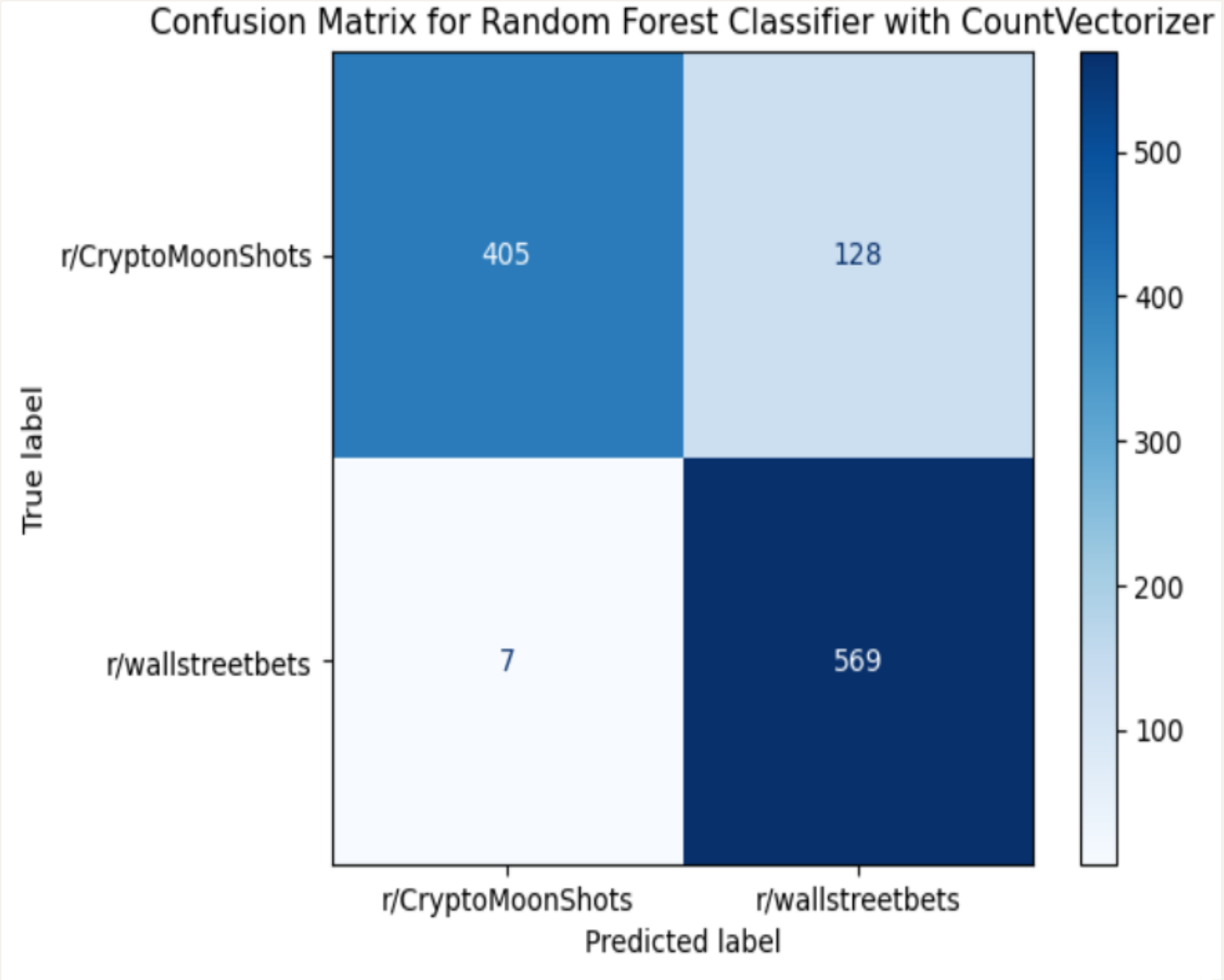
Classification Report for Multinomial Naive Bayes:				
	precision	recall	f1-score	support
0	0.94	0.91	0.92	533
1	0.92	0.94	0.93	576
accuracy			0.93	1109
macro avg	0.93	0.93	0.93	1109
weighted avg	0.93	0.93	0.93	1109



Model 3: Random Forest Tree Classifier with CountVectorizer

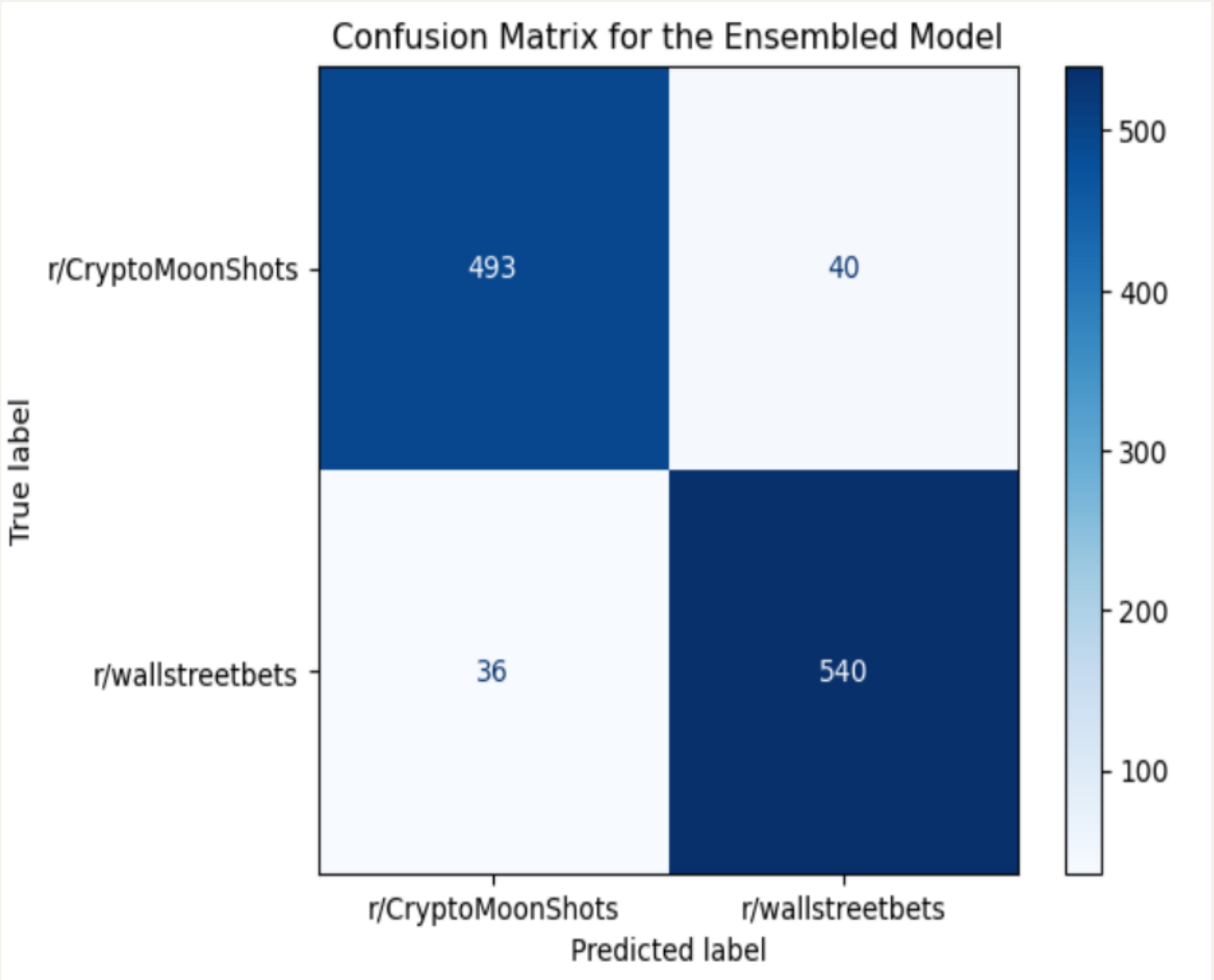
Training Accuracy	Testing Accuracy
0.891	0.878

Classification Report for Random Forest Classifier:				
	precision	recall	f1-score	support
0	0.98	0.76	0.86	533
1	0.82	0.99	0.89	576
accuracy			0.88	1109
macro avg	0.90	0.87	0.88	1109
weighted avg	0.90	0.88	0.88	1109



Model 4: Ensemble Model (Logistic Regression & Multinomial Naïve Bayes)

Testing Accuracy				
0.931				
Classification Report for Ensembled Model:				
	precision	recall	f1-score	support
0	0.93	0.92	0.93	533
1	0.93	0.94	0.93	576
accuracy			0.93	1109
macro avg	0.93	0.93	0.93	1109
weighted avg	0.93	0.93	0.93	1109



Conclusions/Insights

- **Precision, Recall, F1-score, and Accuracy:** The ensemble model has similar performance to the first two individual models, with **high precision of 0.93, recall of around 0.93, F1-score of 0.93**, for both subreddits and an overall **accuracy score of 0.93**.
- This will give more of a balance with the false negatives and positives; with **36 false negatives vs 40 false positives!** Let's **choose this ensemble model** for production!





Thanks!