

WICC x Millennium

Social Forum Discourse on Artificial Intelligence and the Effects on Semiconductor Stocks

Team 3

April 2025



Table of Contents

1. Introduction
2. Exploratory Data Analysis
3. Market Sentiment
4. Combined Analysis
5. Final Takeaways
6. Closing Words

Introduction



Meet the (dream) Team!



**Helen
Bian**

CS '26



**Akshika
Chawade**

CS '26



**Serena
Duncan**

CS '26



**Grace
Jin**

CS + ECE '27



**Esha
Shah**

CS '27



**Grace
Wei**

CS '26



**Natasha
Dilamani**

Our Mentor

Our Project Timeline

Fall 2024

Team formation

Met with Millennium and learned the scope of the project.

Brainstormed ideas.

Jan. - Feb. 2025

Data Collection

Drafted a preliminary investigation plan and received feedback from our mentor Natasha.

Began collecting stock data and online sentiment data.

Feb. - Mar. 2025

Preliminary Analysis

Continued collecting different types of data.

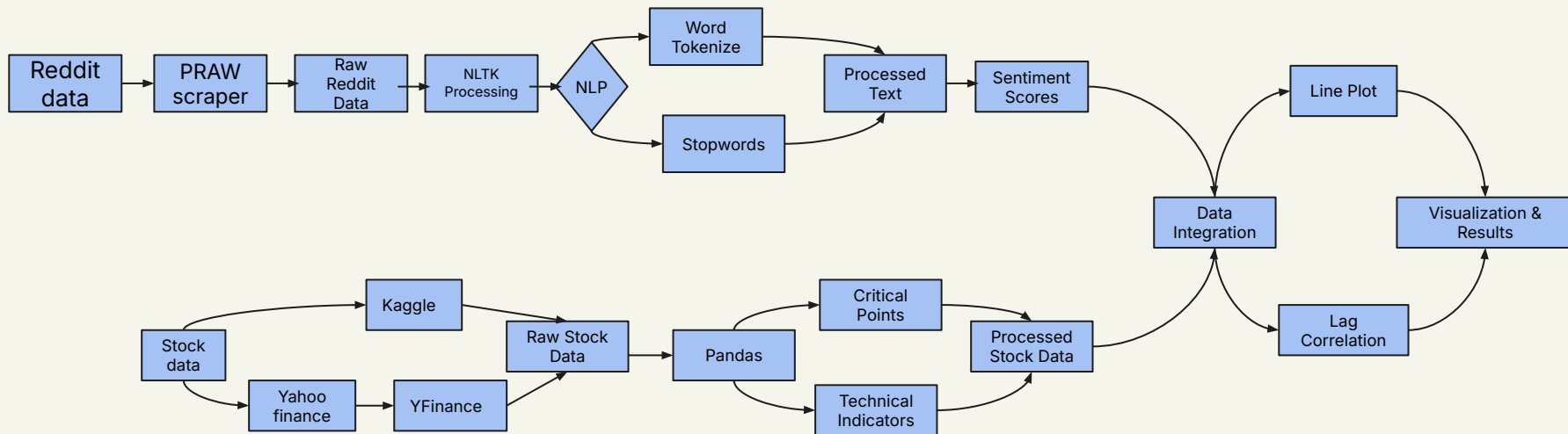
Created programs for sentiment and data analysis.

Apr. 2025

Final Analysis

Collected results from analysis and wrote final conclusions.

Our project architecture



Project Background + Inspiration

- **AI Boom & Market Disruption**

The rapid growth of AI technologies has reshaped the tech industry, spotlighting companies like NVIDIA, essential to AI hardware and software development.

- **NVIDIA's Rise in Popularity**

NVIDIA has become a key player in the AI revolution, with its GPUs driving advancements in machine learning, leading to significant stock growth.

- **Major Market Events**

Notable events, such as the release of **DeepSeek**, caused market disruptions — highlighting the volatility and risks associated with tech stocks.

- **Focus Shift to NVIDIA**

While the initial goal was to explore sentiment effects on the semiconductor industry, the focus narrowed to NVIDIA due to its prominence in recent market movements.

- **Exploring Sentiment & Stock Correlation**

Our project investigates the correlation between Reddit sentiment and NVIDIA's stock price, analyzing how online discussions might influence market behavior.

Goals

Sentiment Correlation

Understand the relationship between sentiment and market behavior

Event Detection

Identify sentiment patterns around major market events

Strategy Insight

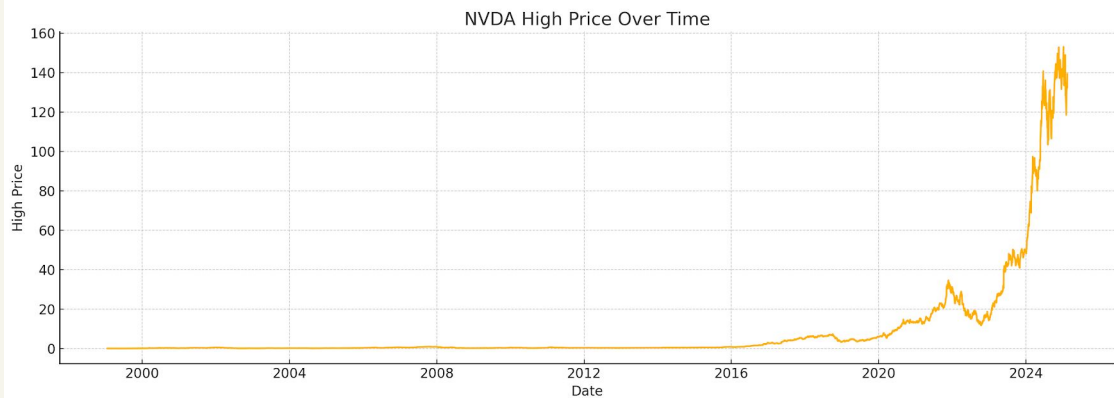
Explore how sentiment trends might inform trading strategies

A Look Into NVIDIA:

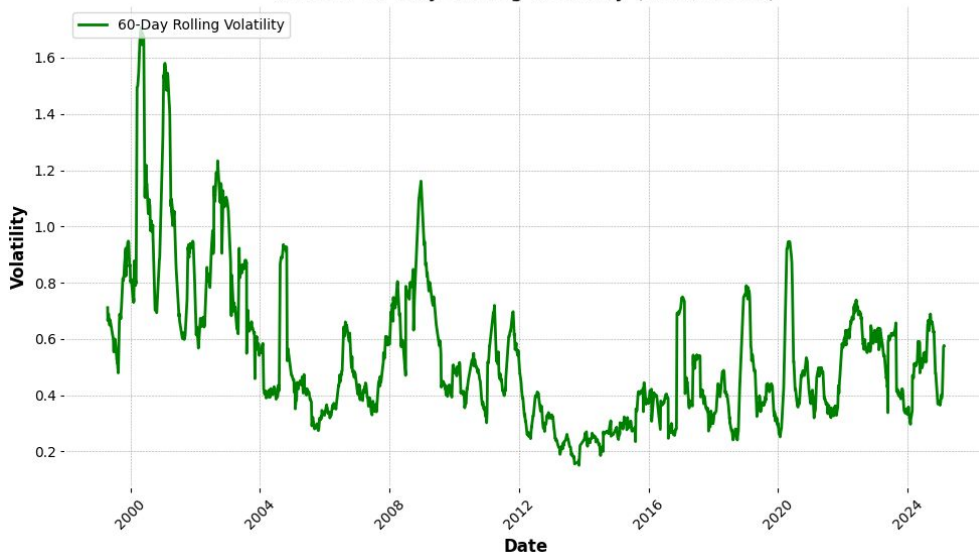
An exploratory data analysis

Observations and Statistics

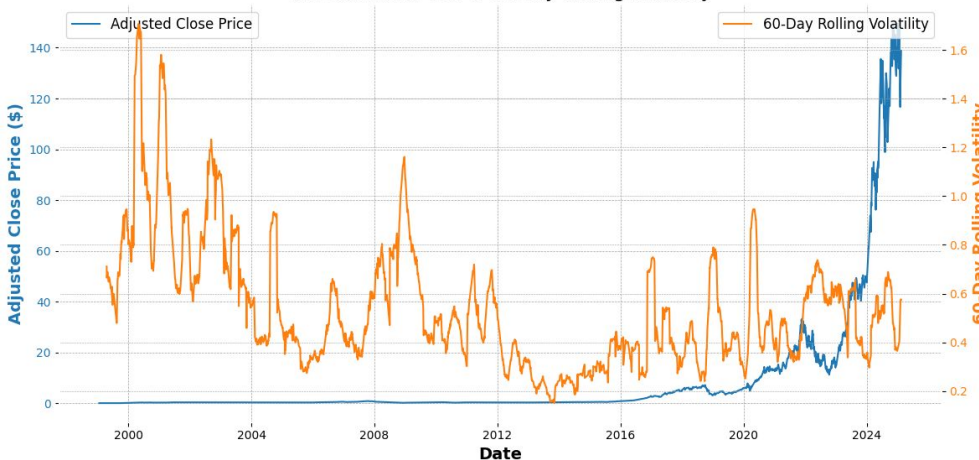
- The stock has experienced explosive recent growth-dramatic upward momentum, particularly from 2020 to 2024, with the most dramatic rise occurring in the 2023-2024 period.
- Price range: From trading below \$20 for many years, the stock has increased to well over \$140 by 2024 (approximately 7-8x increase).
- As the price has increased, the stock has shown greater price volatility, with larger daily and weekly price swings.



NVIDIA 60-Day Rolling Volatility (Annualized)



NVIDIA Stock Price & 60-Day Rolling Volatility



Observations

- Volatility peaks in the early 2000s, but this could be due to the very low market price of the stocks during the time
 - Also coincides with the "Dot-Com Bubble burst" era of the American economy
- Spikes after 2008, 2020
- General downward trend
 - Opposite of upward growth in market price

Narrowing Our Vision: 2020-2025

- AI Boom and the Rise of NVIDIA's Market Influence (2020s)
- More recent and relevant data to understand market behavior
- NVDA stock started dramatically increasing around 2020, which attracted more attention from investors and the public eye, which led for more commentary that we could analyze.
- This time frame encompasses multiple market cycles, including the COVID-19 market crash (March 2020), recovery, which results in diverse data points
- 2020-2025 marks the transition period when NVIDIA shifted from being primarily viewed as a gaming/crypto mining company to becoming the dominant AI computing platform, fundamentally changing investor perceptions.

Market Sentiment



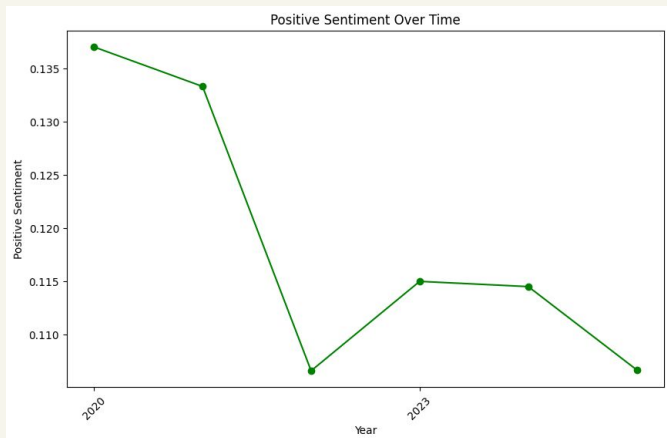
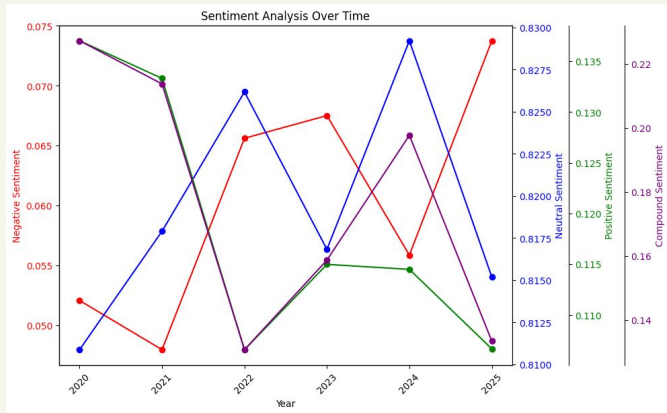
Data Source: Reddit

- Data collected from Reddit posts across subreddits r/stocks and r/investing from the years 2020-2025
- We scraped Reddit with PRAW (Python Reddit API Wrapper) with target keyword filtering ("nvidia","nvda") and post metadata
- We chose Reddit for its less-filtered sentiment that isn't available in traditional financial news sources

```
1 import praw
2 import pandas as pd
3 import datetime
4 import time
5 from tqdm import tqdm
6
7 def scrape_subreddit(subreddit, subreddit_name, keywords, start_date, end_date, limit=500):
8
9     start_timestamp = int(start_date.timestamp())
10    end_timestamp = int(end_date.timestamp())
11
12    subreddit = reddit.subreddit(subreddit_name)
13    posts = []
14
15    query = " OR ".join(keywords)
16    print(f"Searching for: {query}")
17
18    search_count = 0
19
20    for post in tqdm(subreddit.search(query, sort='new', time_filter='all', limit=None),
21                    desc=f"Searching posts", unit="post"):
22        search_count += 1
23
24        if len(posts) >= limit:
25            break
26
27        if start_timestamp <= post.created_utc <= end_timestamp:
28            title_lower = post.title.lower()
29            selftext_lower = post.selftext.lower() if hasattr(post, 'selftext') else ""
30            matched_keywords = []
31            for keyword in keywords:
32                keyword_lower = keyword.lower()
33                if keyword_lower in title_lower or keyword_lower in selftext_lower:
34                    matched_keywords.append(keyword)
35
36            if matched_keywords:
37                posts.append({
38                    'post_id': post.id,
39                    'title': post.title,
40                    'author': str(post.author) if post.author else "[deleted]",
41                    'created_utc': post.created_utc,
42                    'date': datetime.datetime.fromtimestamp(post.created_utc).strftime('%Y-%m-%d %H:%M:%S'),
43                    'score': post.score,
44                    'num_comments': post.num_comments,
45                    'permalink': f"https://www.reddit.com/{post.permalink}",
46                    'selftext': post.selftext if hasattr(post, 'selftext') else "",
47                    'matched_keywords': ", ".join(matched_keywords)
48                })
49            print(f"Found matching post from {datetime.datetime.fromtimestamp(post.created_utc).strftime('%Y-%m-%d %H:%M:%S')}")
50
51    if search_count % 10 == 0:
52        time.sleep(1) # sleep to avoid hitting rate limits
53
54    print(f"Collected {len(posts)} posts matching the criteria")
```

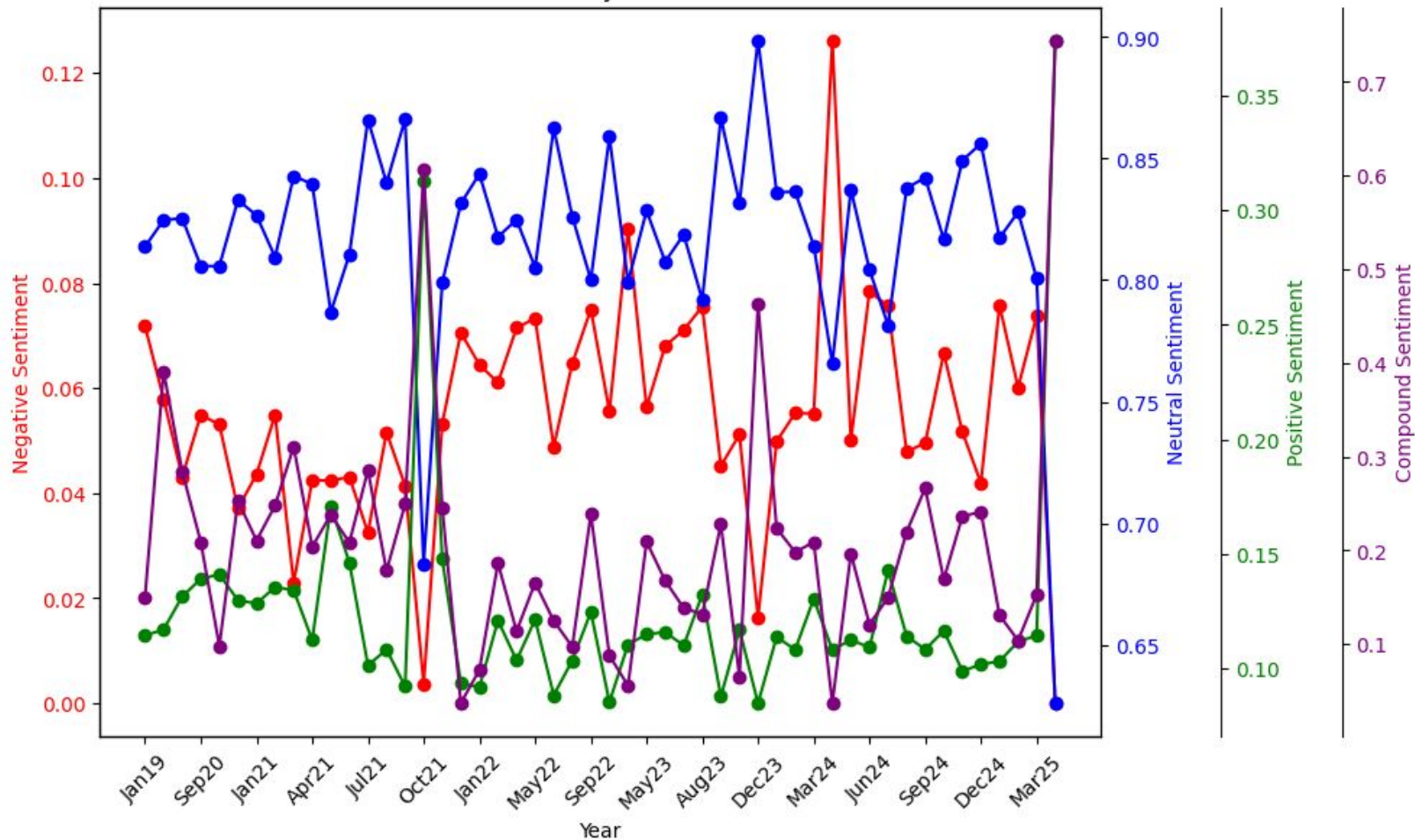
Critical Part in Data Collection
Script

Sentiment Over Time

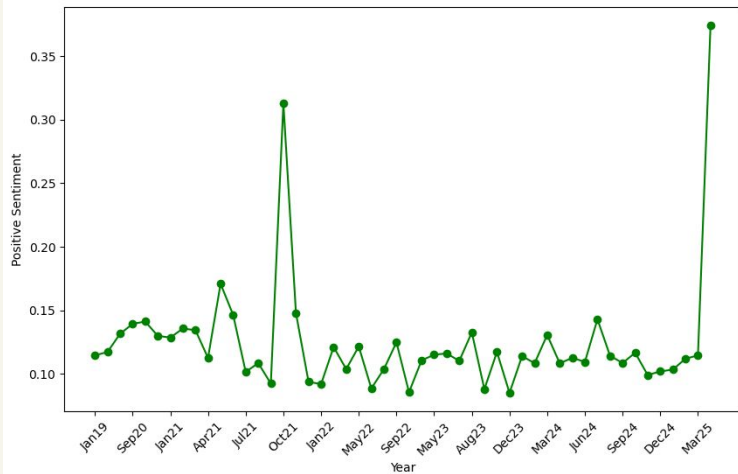


- Used NLTK to do the following (2021-2025):
 - Find average sentiment across years
 - Find top 15 words across years
 - Find average sentiment per word
- Why NLTK?
 - Easy-to-use sentiment analysis tool
 - Built-in tools for:
 - Stopword removal
 - Tokenization
 - Sentiment scoring
 - *Limitation*: Basic sentiment analysis - might miss context or sarcasm
- What did we find?
 - Positive sentiment around NVIDIA highest in 2020
 - Drops **significantly** from 2020-2022
 - Slowly recovers through 2023-2024 but never reached 2020 levels again
 - Drops again in 2025

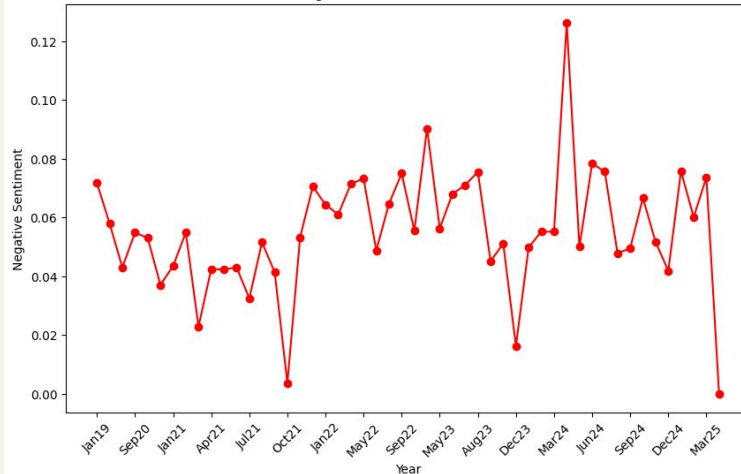
Sentiment Analysis Over Time



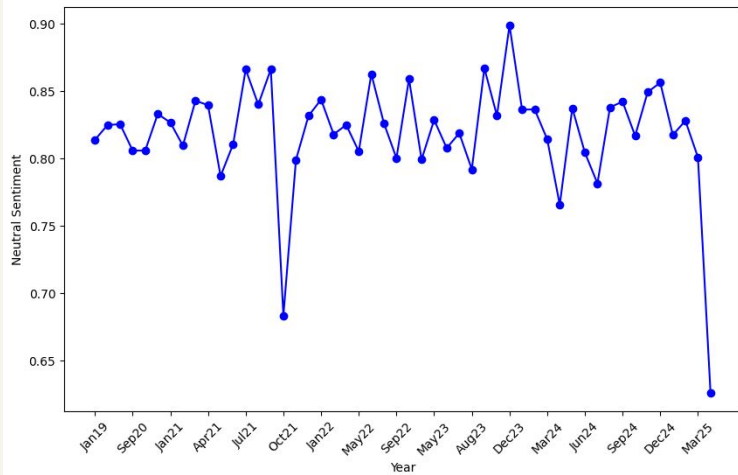
Positive Sentiment Over Time



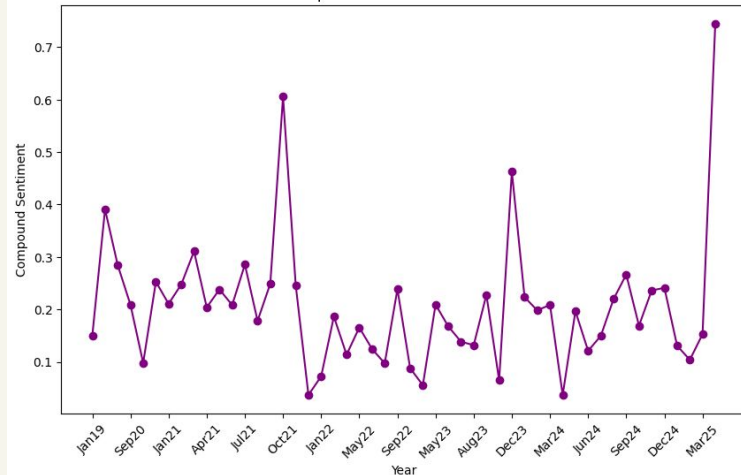
Negative Sentiment Over Time



Neutral Sentiment Over Time



Compound Sentiment Over Time



Prevalence of AI

- What we looked at:
 - Top 15 most common words in r/stock or r/investing reddit posts & comments about NVIDIA (2021 - 2025)
- Key trends we saw:
 - **Early years (2021-2022)** → Conversations dominated by:
 - *Stock splits, shares, crypto, gaming, mining*
 - **Starting in 2023** → Clear rise of AI mentions
 - "AI" enters top words in 2023
 - Becomes **2nd most common** word by 2024
 - **Most common** word by 2025
- Other Emerging Keywords:
 - *Deepseek* (AI-related) appears in 2025
 - Consistent growth of *NVIDIA* mentions over time (likely tied to AI leadership)

```

[5] def process_text_data(df, text_columns):
    for i, row in df.iterrows():
        combined_text = ''
        for col in text_columns:
            combined_text += ' ' + str(row[col]) for col in text_columns
        month = datetime.utcfromtimestamp(row['created_utc']).strftime('%b%Y')
        sentiment = sia.polarity_scores(combined_text)

        tokens = word_tokenize(combined_text.lower())
        filtered_tokens = [word for word in tokens if word not in stop_words and word not in punct and word.isalpha()]

        word_freq = Counter(filtered_tokens)
        top_15 = word_freq.most_common(15)

        if month not in month_sentiment_sum:
            month_sentiment_sum[month] = sentiment.copy()
            month_sentiment_count[month] = 1
        else:
            for k in sentiment:
                month_sentiment_sum[month][k] += sentiment[k]
                month_sentiment_count[month] += 1

        for word, count in top_15:
            if month not in word_counts:
                word_counts[month] = {}

            if word in word_counts[month]:
                word_counts[month][word] += count
            else:
                word_counts[month][word] = count

            if word in word_to_sentiment:
                for k in sentiment:
                    word_to_sentiment[word][k] += sentiment[k]
                    word_to_sentiment[word]['count'] += 1
            else:
                word_to_sentiment[word] = sentiment.copy()
                word_to_sentiment[word]['count'] = 1

[6] posts_df = pd.read_csv('investing_posts_nvidia_2020-2025.csv')
comments_df = pd.read_csv('investing_comments_nvidia_2020-2025.csv')

# posts_df = pd.read_csv('stocks_posts_nvidia_2020-2025.csv')
# comments_df = pd.read_csv('stocks_comments_nvidia_2020-2025.csv')

process_text_data(posts_df, ['title', 'selftext'])
process_text_data(comments_df, ['body'])

[7] for word, sentiment_data in word_to_sentiment.items():
    count = sentiment_data.pop('count')
    word_to_sentiment[word] = {k: sentiment_data[k] / count for k in sentiment_data}

[8] month_to_sentiment = {
    month: {k: month_sentiment_sum[month][k] / month_sentiment_count[month] for k in month_sentiment_sum[month]}
    for month in month_sentiment_sum
}

[9] for month in word_counts:
    month_to_word[month] = sorted(word_counts[month].items(), key=lambda x: x[1], reverse=True)[:15]

[10] print(month_to_sentiment)
print(word_to_sentiment)
print(month_to_word)

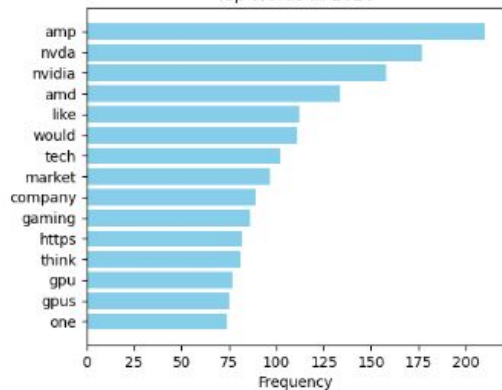
[11] # {'Apr25': {'neg': 0.8522, 'neu': 0.8466, 'pos': 0.10180333333333334, 'compound': 0.36876333333333334, 'Mr25': {'neg': 0.859831, 'advice': {'neg': 0.8345454545454546, 'neu': 0.8485454545454547, 'pos': 0.12518181818181817, 'compound': 0.7268181818181818}, 'Apr25': {'neg': 0.8345454545454546, 'neu': 0.8485454545454547, 'pos': 0.12518181818181817, 'compound': 0.7268181818181818}, 'Apr25': {'advice': 3, 'investments': 4, 'etfs': 3, 'upst': 2, 'etly': 2, 'hrt': 2, 'help': 1, 'investment': 2, 'struggling': 1}}

[12] normalized_word_counts = {
    month: word_count / sum(month_word_counts.values()) for month, month_word_counts in word_counts.items()
}

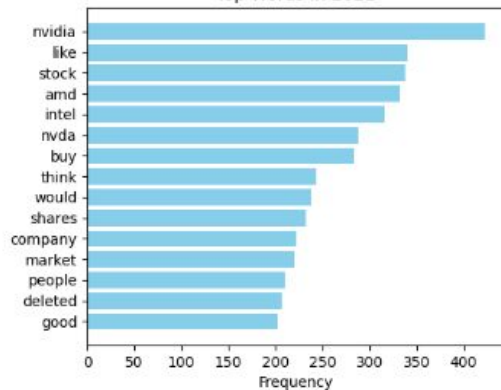
```

Critical Part in Sentiment Analysis Script

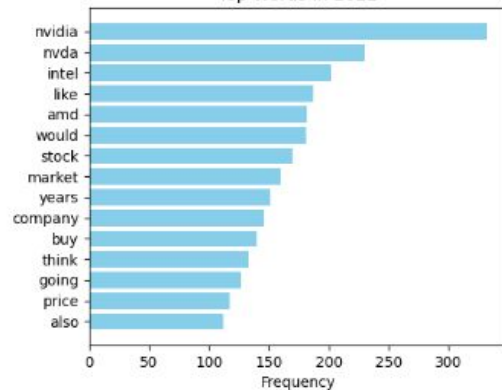
Top Words in 2020



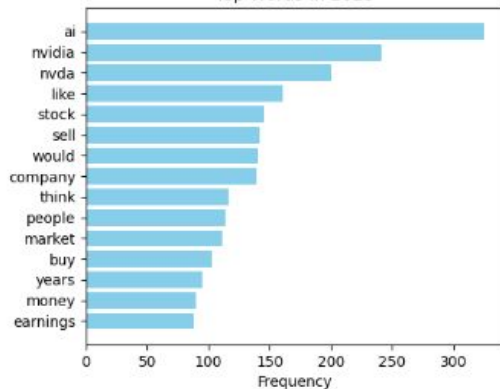
Top Words in 2021



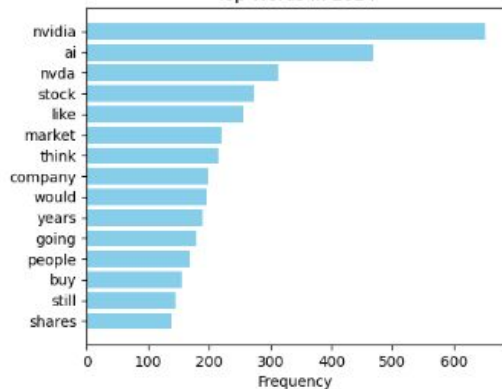
Top Words in 2022



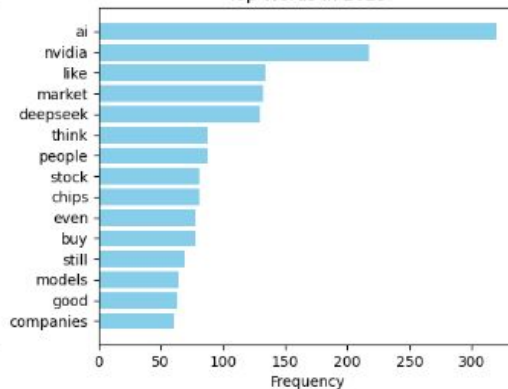
Top Words in 2023



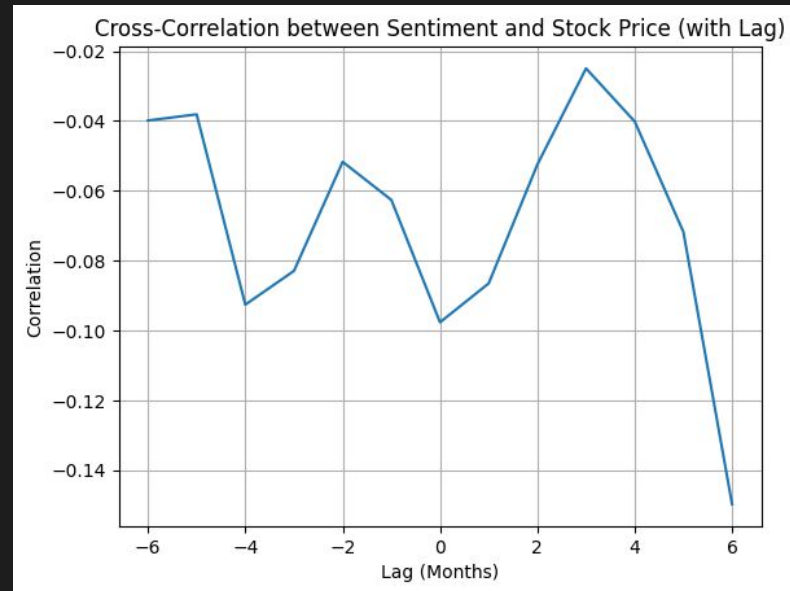
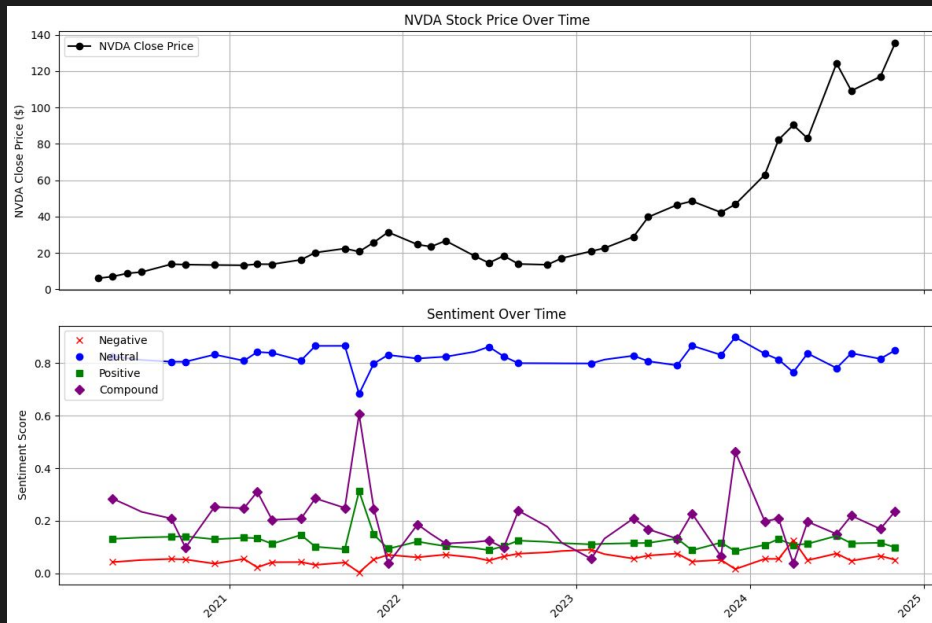
Top Words in 2024

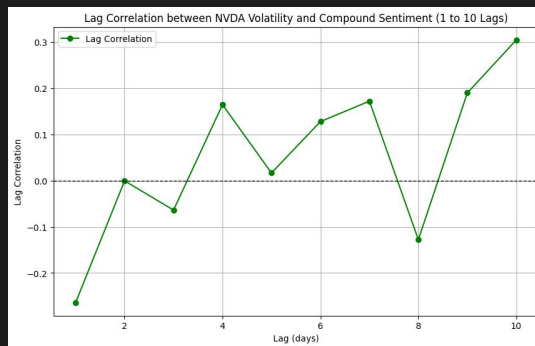
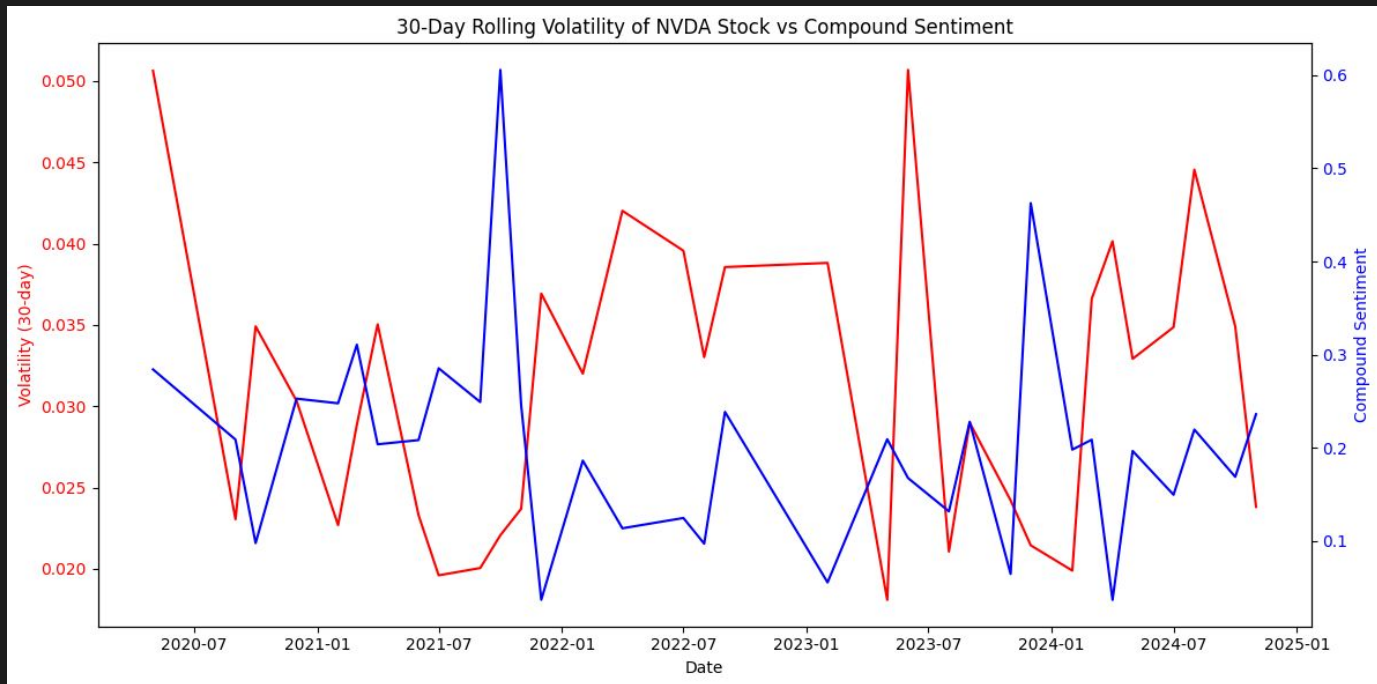


Top Words in 2025



Combined Analysis





Summarizing Findings

Sentiment does not significantly affect stock growth

The correlation between sentiment (both positive and negative) and NVDA stock price growth is minimal. Although fluctuations in sentiment were observed, they did not show a clear, direct impact on the stock's short-term or long-term performance.

Sentiment and volatility correlation is weak

While there are periods of heightened sentiment, especially with terms like "AI," the sentiment does not show a strong correlation with stock price volatility. The stock price volatility does not seem to increase during periods of high sentiment, suggesting that sentiment alone is not a key driver of price fluctuations for NVDA.

Limitations

- Difficult to collect comprehensive and balanced data set (e.g. samples across years)
- Issues collecting data with API, and having to pivot at the last minute → difficult to collect many samples
- Uses sentiment data from a single platform, which may not fully capture broader market sentiment
- Difficult distinguishing between social media sarcasm and genuine analysis without more advanced NLP techniques

Beyond Our Project

To improve our current project, we would:

- Diversify our data sources beyond Reddit (e.g. Twitter, YouTube comments, paywalled stock forums)
- Use fine-tuned transformer models (e.g., BERT, RoBERTa) for more nuanced sentiment classification
- Examine the effects on other companies to gain a better understanding of the semiconductor and GPU manufacturing market, as well as the tech industry beyond computer hardware

Final Takeaways

- Learned a lot of new skills — got a lot of exposure to new technologies
 - Got experience web scraping, using external APIs, NLP and text processing techniques
- Learned better data wrangling techniques
- Learned about the stock market and visualizations and metrics for analyzing stock market events

Thank you!