**Abstract**

In the 2020 US Election there was a huge amount of discussion on Twitter, much of it hostile and much of it based on conspiracy. I will use Natural Language Processing, and specifically Sentiment Analysis to interrogate and analyse a large corpus of Twitter data from the period.

I first had to distil the large Twitter dataset into parts that were easier to process, split by prominent figures (Trump, Biden, Pence, Harris). These data frames were then run through various pipelines to explore the Sentiment Analysis and Conspiracy Prevalence.

The analysis showed that Trump as a figure trended towards negative sentiment and conspiracy prevalence, going so far as to increase the negative sentiment and conspiracy prevalence of those he was paired with. Additionally, retweets were found to have a major effect on sentiment and conspiracies.

There were various questions answered, primarily concerning the negative influence of Trump. Some areas could be improved, specifically the accuracy of the Sentiment Analysis tools and a more accurate means to determine tweets with specific conspiracy theories.

# Contents

## Table of Figures

## Table of Tables

## 1.1 Introduction

Social media as technology has grown since its inception in the 90s with the widespread adoption of the World Wide Web. From the beginnings of relatively simple sites such as SixDegrees.com and GeoCities, monoliths of the industry have emerged such as Facebook and Twitter. From each company's Second Quarter 2022 Results, Twitter reports 237.8 million daily users (Twitter 2022) while Facebook reports a staggering 1.97 billion daily users (Facebook 2022).

While these platforms have great benefits such as social connectivity, social involvement, information attainment, and entertainment, there are also many risks associated with their use (Khan et al 2014). These risks can include overuse, mental health issues, social problems, and privacy (ibid.). This study will focus on the societal drawbacks that can be enabled through the use of social media, using a set of Twitter data gathered from March to July of 2020, exploring the US election of that year and the prevalence of conspiracy theories in that period.

In the 2020 election, there was a huge amount of disinformation being propagated on social media, this was in no small part due to the controversial figure of the incumbent American President Trump who promoted much of the disinformation and conspiracy theories.

Using powerful tools included with the Natural Language ToolKit (NLTK) for Natural Language Processing (NLP), I hope to explore the prevalence of specific conspiracy theories during the election in relation to specific candidates or pairings of candidates, as well as any trends in sentiment analysis with these candidates or pairings of candidates.

## 1.2 Aims and Objectives

- Develop a tool for curating a large corpus of tweets and turning it into a manageable data set that fits my requirements.
- Explore which figures and pairings of figures (Trump, Biden, Pence, Harris) are most commonly associated with conspiracy theories/certain sentiment analyses.
- Evaluation of the VADER sentiment analysis tool for interpreting conspiracy theories on Twitter.

# 2 Background
## 2.1 Previous Work
### 2.1.1 Conspiracy Theories

There are three main conspiracy theories explored in this project; covid-related, QAnon, and Ukraine & Russia.

1. Conspiracy theories related to covid-19 are very broad and come in several varieties. There are those who believe covid is a bioweapon engineered in China, those who think the elite are using the opportunity to perform unnecessary vaccinations for some nefarious purpose, those who think the whole pandemic is a lie and that the disease doesn't actually exist or is seriously exaggerated. While not strictly a conspiracy theory while President, Trump contributed to spreading a great deal of misinformation including the use of dubious medical treatments (Mahase 2020. Additionally, the 2020 election was taking place at the height of the first wave of the pandemic.

   The wanton spread of this conspiracy theory throughout social media and society at large resulted in damage to public health, as certain groups of people would not engage in social distancing or other preventative measures and later refused to receive the vaccine (Romer & Jamieson 2020).

2. QAnon is a conspiracy theory that emerged on a message board and describes a wide-reaching conspiracy being perpetrated by the elite, in which Donald Trump is a sort of saviour. The theory is too complicated and bizarre to explain succinctly but with its close ties to Trump it has an important place in this project. Outside the scope of the data used in this project, QAnon had a huge impact in the real world as the so-called 'qarmy' were heavily involved in the January 6[th] Insurrection in 2021.

   Perhaps more than any other conspiracy theory QAnon has benefitted from the spread of social media. Not only was it born on the image board 4chan, but from its inception to June 17 2020 the Institute for Strategic Dialogue recorded "69,475,451 million tweets, 487,310 Facebook posts, and 281,554 Instagram posts mentioning QAnon-related hashtags and phrases" (ISD 2020).

3.  The third conspiracy theory I have chosen to examine is regarding Joe Biden and the allegations that he had used his position as Vice-President to pressure Ukraine into dropping a corruption investigation that would have put his son Hunter Biden under scrutiny. This was pushed as an attempt to discredit Biden's presidential bid.
    While not

## 2.1.2 Conspiracy Theories on Social Media

Conspiracy theories as a whole are a complex thing to define, with the task not growing any easier when trying to pin down a specific theory. One common thread of most conspiracy theories is the belief in a truth that does not align with the mainstream and often factual truth, often with the belief that some other party is suppressing the theorist's 'truth' for some nefarious reason (Marcellino et al 2021).

This leads to an entrenching of their position as the believers reject information that does not align with their theories and seek out sources that reinforce their beliefs. This process has been made much more prevalent with the use of social media. As a user of social media can curate who they follow, echo chambers can emerge. An echo chamber is a community where dissenting voices are not heard and the members' beliefs can be reinforced (Wardle & Derakhshan 2017).

Social media is a medium for which a great deal of misinformation can be spread, with people who use social media as a source of news increasing the chance that they will hold beliefs in a conspiracy theory (Stecula & Pickup 2021).

There are many ways that conspiracy theories can be spread across social media:

- Bots can be used by those within the conspiracy community or outside agents who wish to fan the flames of that theory (Shao et al 2018)
- Organised raids can be used to artificially inflate the reach of rumours or misinformation with reposting, resharing, or retweeting as well as general public discussion (Wardle & Derakhshan 2017)
- Hashtags can be forced to trend by concerted efforts by a large community, or an existing hashtag can be hijacked by posting irrelevant or harmful information using an already popular hashtag, with the intention of diluting the

content of the hashtag or with malicious intent for those participating in the conversation (Marwick & Lewis 2017)

### 2.1.3 Natural Language Processing

Natural Language Processing (NLP) is a wide-ranging computational technique and research domain, where the goal is for the comprehension of natural human language by a machine along with performing some level of linguistic analysis (Liddy 2010). This can be attempted and achieved with varying levels of success through a range of techniques.

A machine capable of understanding the idiosyncrasies of human speech is a tall task, as there is a wide range of quirks that humans simply take for granted, such as synonyms, idioms, and humour.

## 2.2 Python Libraries

As this project was written entirely in Python using Jupyter Notebooks, all of the modules and libraries utilised are Python modules and libraries.

### 2.2.1 File Management

os - This module allows for the use of operating system functionality. This project makes use of the listdir function to assist in opening all the files in a given directory.

gzip - This module allows for the compression and decompression of gzip files. This project makes use of the open function to read each individual tweet stored within the many gzip files provided.

### 2.2.4 Natural Language Processing

NLTK - The Natural Language Tool Kit is a suite of libraries for NLP. It is suitable for those who are first experimenting with NLP as well as practitioners and those conducting research into NLP (Chowdhary 2020), so it is no doubt ideal for a first project in the field.

An incredibly important suite, it includes all the base tools to perform NLP tasks that are useful to this project such as tokenization, removal of stop words, and with the

VADER model a great sentiment analysis tool that sits at the heart of much of this project. (Bird, Klein & Loper 2009)

<u>VADER</u> - Vader is a sentiment analysis tool that takes a lexicon and rules-based approach, specifically designed for use with social media posts. Additionally, VADER uses intensity scores (ranking just how positive or negative something is) rather than just a binary (deciding if something is either good or bad), which allows for a more nuanced result for sentiment analysis (Hutto & Gilbert 2014).

As sentiment analysis is at the very heart of this project, it would have been hard to get anything completed without an NLP tool of some kind and VADER was the one that fit the job best. Using the 'polarity_scores' method on a string returns a dictionary containing four scores; positive, neutral, negative & compound (a combination of the first three).

Though it stretched beyond the scope of this project, the tool could have been adapted to perform even better by adding additional conspiracy-related words with matching sentiment scores to the VADER lexicon.

## 2.2.3 Data Manipulation & Visualisation

<u>pandas</u> - pandas is a very powerful library allowing for data manipulation and analysis. One of the most essential libraries for this project, as it was used to create all of the data frames for the project (pandas 2020, McKinney 2010).

Using simple and powerful tools within pandas, the data could then be manipulated and used in concert with Matplotlib to generate data visualisations.

Additionally using the 'to_pickle' method, the many large data frames that were pulled from the misinformation callout corpus were able to be saved, so the extraction only needed to be performed once.

<u>Matplotlib</u> - Matplotlib is a library for the creation of a wide range of data visualisations. While the charts created in this project were all relatively rudimentary, the library is a powerful tool with the potential for creating animated and interactive visualisations (Hunter 2007).

As the project involved analysing millions of tweets, it would have been a great hindrance to do so without a good means of displaying the data in a more visual medium.

### 2.2.7 Tweepy

Tweepy is an easy-to-use library to access the Twitter API and gather the tweets that Twitter has made publicly available (Roesslein 2020).

As I was only engaging in a small amount of gathering of my own Twitter data, Tweepy was a perfect fit as I did not have to learn to use the Twitter API for a relatively small part of the overall project.

# 3 Data Collection and Organisation

## 3.1 Existing Corpus

The majority of data that is used in this study was taken from a huge dataset of Tweets ordered by month and year, ranging from 2019 to 2022. Within each month was a json.gz file, a compressed json file that contained millions of tweets. The data was gathered searching for tweets for which the text included any of a list of search terms relating to disinformation, misinformation, and conspiracy theories amongst other things. This was done using the Sentinel platform (Preece et al 2017).

Due to this expansive source of Twitter data, the necessity for my own data collection was quite minimal. However, for Part 5 there was some use of the Twitter API to collect tweets from August-September 2022, as the misinformation callout corpus did not include tweets from that period.

The misinformation callout corpus comprised of tweets gathered every day in the months of March, April, May & June 2020. The June set was missing entries from 15th-28th, so I decided to simply drop the last few days of June as well, meaning I used data stretching over a total of 15 weeks.

All the json files in this period contained over 8 trillion tweets, so a careful selection would have to be made. This was not only due to the prohibitively large amount of data, but also the fact that the vast majority of it would have been irrelevant to the subject at hand.

## 3.2 Curating the Corpus

### 3.2.1 Dividing the Corpus by Important Figures

I decided to focus on important figures in the 2020 election, selecting both presidential candidates, Trump and Biden, and their running mates, Pence and Harris. From the main datasets tweets would be selected that mentioned either one or two of these figures, and then organised into weekly datasets for easier management.

The option of examining tweets that mentioned either three of the figures or all four together was explored, but the sizes of those curated datasets were all too small to be worth exploring further. Additionally, when I first gathered my own data frames

from the original corpus, I organised the data by month but was later changed to be organised by week to be able to more carefully examine the data.

As can be seen in Figure 1 and Figure 2 when organised by month, outlying data points would obfuscate a quarter of the data. While organised by week it gave a more precise impression of what was occurring, with the outlier only affecting a single week.
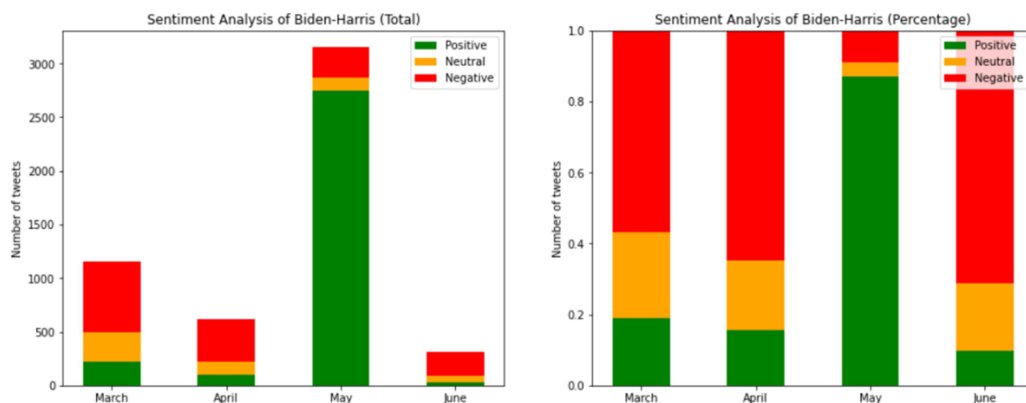


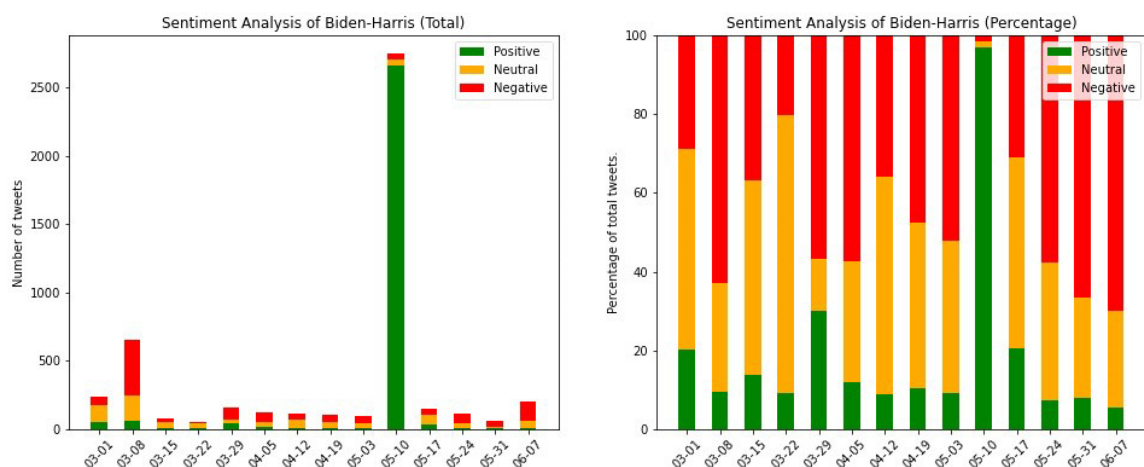Figure 1. General Sentiment Analysis of Biden-Harris (by month)



Figure 2. General Sentiment Analysis of Biden-Harris (by week)

Once each of the json.gz was unzipped each of these files ranged in sizes from 2-5 gigabytes, which was far too large to allow unzipping the entire dataset and working from that. Even unzipping the files by month proved to be a cumbersome and untenable process.

Instead, I iterated through the directories the files were saved to and used the gzip library 'open' method to read each file within each directory without saving it. While this was still time-consuming, it saved a huge amount of disk space.

### 3.2.2 Columns for the Data Frame

The corpus had a huge range of potential data points for every tweet contained within, so for my own datasets I chose to focus on the following fields, some of which were taken directly from the corpus and some of which were derived from data within it;

1. user_id: the id of the individual user who made the tweet. Note that the intention of this was to check for 'super users' whose tweets are prevalent in the dataset, and not as an attempt to identify the anonymous users.
2. tweet_id: the id of the individual tweet.
3. text: the content of the tweet. From this field, the names, hashtags, sentiment, conspiracies, trumpTags & bidenTags were derived.
   Certain common pre-processing techniques such as removing punctuation and lower casing were not performed, as these elements are important parts of the VADER sentiment analysis process.
   Mentions (any token beginning with '@') and URLs (any token beginning with 'http') were removed, as they contributed
4. date: the date and time the tweet was made
5. names: which of the important figures were mentioned in this tweet
6. hashtags: any hashtags in the tweet, ignoring individual '#'s
7. sentimentA: the compound sentiment score as provided by VADER, converted to either 'Positive' (if score is >= 0.5), 'Neutral' (if score is > -0.5 and < 0.5) or 'Negative' (if score is <=-0.5
8. sentimentB: an alternate form of recording the sentiment. Still using VADER, but instead of calculating based on the compound score, simply recording whether the positive score or negative score was higher.
9. retweet_count: the number of retweets the tweet received, if any
10. conspiracies: if any hashtags that were associated with a specific conspiracy theory were mentioned in the body of the tweet, then the name of that conspiracy theory was added to this list

11. trumpTags: an integer showing how many Trump hashtags were mentioned in the tweet
12. bidenTags: an integer showing how many Biden hashtags were mentioned in the tweet

Deciding on these columns was an iterative process, as the focus and scope of the project expanded more columns were added. For instance, initially the data that appears in the final 'conspiracies' column was being pulled from hashtags while performing NLP in later sections. Moving the generation of the hashtags and making it a permanent column of the data frame made it easier to explore the data later on.

### 3.2.3 Unused Fields
There were some fields included in the data frames which I had intended to utilise, but in the end did not have enough time to implement them in a meaningful way.

- User ID was going to be used to find if there were any 'super users' who posted more regularly and in doing so skewed the data, and see if 'super users' were more common in a specific figure/pairing or conspiracy theory. However more things proved more informative and this angle was dropped.
- TrumpTags and BidenTags were the first steps in trying to organise tweets by Trump supporters and Biden supporters, in a similar way to how conspiracy theories were organised by hashtag. Unfortunately, this approach was far too rudimentary, and to properly realise this idea would have required a heavy investment of time and machine learning.
- SentimentB was intended to be used with the General Sentiment Analysis to show sentiment distributions of Positive/Negative instead of Positive/Neutral/Negative. While I did create these graphs, I did not find that they added a significant amount of further analysis so did not include them in this report.

## 3.3 Categorising Hashtags
### 3.3.1 Hashtag Generation
To create the 'conspiracies' field for my dataset I needed some way to identify which tweets were talking about specific conspiracy theories.

As I was working with such a large dataset, I chose a relatively simple solution. As I already planned to gather the hashtags that were present in each tweet I decided to utilise that data, as hashtags are good indicators of a tweet's overall message.

n.b. I also utilised a similar workflow to populate the bidenTags and trumpTags fields, but as they ended up being unused in the final project I won't describe it in detail.

### 3.3.2 Choosing Conspiracy Theories

I chose some theories which I thought may have been prevalent during this period, with an initial selection of; covid, QAnon, Ukraine (the Hunter-Biden scandal), Russia (Russian interference in the election), fraud (voter fraud), and climate change denial.

Through performing the following steps a few things became clear that led to a refinement of these topics.

1. Many of the hashtags used in the Ukraine and Russia conspiracy hashtags were identical, so I combined them into one topic.
2. While there were large accusations of voter fraud during the 2020 election that did not begin until the time surrounding Election Day itself, and that was over four months after the end of the available data. Due to this, the number of relevant tweets was rather low and I removed them.
3. While climate change denial is endemic on Twitter, it was not a major topic of discussion about the 2020 election when compared to other conspiracies and so I also removed climate change denial from the running.

This left covid conspiracies, which were rampant as it was still the height of the pandemic; QAnon conspiracies, which were similarly popular due to their intrinsic link to Donald Trump; and the renamed ukr_ru category which was popular due to the linked event's impacts on the 2020 election.

### 3.3.3 False Positives

Hashtags are by no means a perfect method of categorising tweets. While some hashtags are used by conspiracy theorists, they are just as likely to not be. Something like '#maga' fits in this category, as while it is commonly used within the QAnon conspiracy theory, it is also used widely by Trump supporters as a whole,

and not all tweets supporting Trump are also discussing QAnon. Additionally, somebody could be using a hashtag that they do not agree with if the tweet's intention is derision or if they are attempting to 'hijack' a hashtag.

Due to this, I was careful to only select hashtags which had a strong bias in them implicitly. For instance '#chinavirus', '#ccpvirus', and '#chinaliedpeopledied' all indicate a belief that China was behind the covid-19 pandemic. Additionally, there were a great number of strange, almost coded hashtags which were in wide use by members of the QAnon conspiracy theory such as '#thestorm','#wwg1wga', and'#thegreatawakening'.

Targeting such pointed or coded tweets such as these reduced the likelihood that somebody was merely discussing the conspiracy theory. These conclusions were drawn by using concordance to check the use of hashtags in context within their tweets.

### 3.3.4 Hashtag Extraction Methodology

To generate the final lists of conspiracy hashtags I created initial lists from two sources;

1. Looking at the most popular overall hashtags from 'curated dataset a' and taking any hashtags that were clearly related to a specific conspiracy theory, backed up by exploring them with concordance.
2. Finding lists of hashtags that had been compiled in similar research which listed hashtags commonly used by the specific conspiracy theories I followed (Dilley, Welna & Foster 2021).

With these initial lists, I then ran the entire curated dataset through a workflow that created a dictionary for each conspiracy theory recording the number of occurrences of any hashtags that appeared alongside my initial list.

This new list required a bit of curating as they often included hashtags that were either not intrinsically linked to the conspiracy theory or were not relevant at all. To accomplish this curation I took hashtags from the new list that I wanted to check and looked up any occurrences of the hashtag in the corpus using concordance from NLTK.

I then combined the initial list with the new list and ran the combined list through the whole process again. I only had to perform this process twice before no new hashtags were found.



Figure 3. Hashtag Extraction Workflow

## 3.4 Corpus Results

Once I ran the large dataset through this workflow, I was left with a total of 7,445,214 tweets, distributed between the individual and paired figures as follows;

| Trump | Trump-Biden | Trump-Pence | Trump-Harris | Biden |
|---|---|---|---|---|
| 6,441,479 | 276,013 | 32,789 | 9,953 | 491,877 |
| Biden-Harris | Biden-Pence | Pence | Pence-Harris | Harris |
| 5,121 | 503 | 26,882 | 32 | 34,709 |

Table 1. Size of figure/pairing data frames

Due to the low number of tweets in the Biden-Pence & Pence-Harris datasets, I chose to remove them, as they would be unlikely to show any interesting results.

## 3.5 Gathering 2022 Data

On August 8[th] 2022 Trump's Mar-a-Lago estate was raided by the FBI in search of classified documents that had allegedly been illegally retained after he failed to be re-elected for a second term. This caused an outcry from Trump supporters, so it seemed a relevant topic to explore and compare the state of tweets currently to the tweets during the 2020 elections.

In the misinformation callout dataset there had been tweets gathered for June & July 2022, but nothing beyond those months, so I chose to gather my own data using the Twitter API and Tweepy library.

## 3.5.1 Twitter API

The Twitter API allows access to all extant public tweets and a great deal of data associated with each Tweet. I applied for the Twitter Developer Account and was granted the necessary API keys to begin making requests to the API and gathering data. As I only had the Essential Access Level available to me I was limited to a maximum of 500,000 tweets per month and only being able to access tweets from the previous 7 days.

Due to this the data I personally gathered was not as extensive as the misinformation callout dataset. I managed to gather tweets from August 5th-September 1st with 25,000 tweets per day.

## 3.5.2 Tweepy

I used the Tweepy library for easy access to the Twitter API. To try and keep my data as similar to the other dataset as possible I used as many of the same search terms as I could. Due to the Essential Access Level, there were limits on the length of the search query string and I had to truncate the list. The only missing terms however were related to Europe, e.g. 'european False', 'european Fraud', 'european Hoax'.

As I would be further curating these tweets to split the gathered corpus into the same important figures that the 2020 dataset is split into, the removal of these European tags seemed inconsequential as the chances of them being related to any of the US figures were low.

Due to user error and misunderstanding of the various fields available to be pulled from the Twitter API, this set of data does not record how many retweets a tweet received, or if it was a retweet at all.

### 3.5.3 Most Popular Times

While it did not make much of a difference given the relatively small amount of data I was able to gather, I set out to determine the best time of day to gather tweets. This was simply done by iterating through all the data frames and incrementing the correct hour in an array.

The resulting graphs showed the most popular times as between 10 pm-3 am UTC, which would be 5 pm-10 pm EST or 2 pm-7 pm PST.

**Figure 4. Popular Tweeting Times**

# 4 2020 Election Dataframes

This data was explored in many ways and went through many changes throughout the project, but can be broadly described in the following ways;

1. General Sentiment Analysis
2. Hashtag Popularity
3. Conspiracy Theory Prevalence
4. Frequency Distribution
5. Concordance

## 4.1.1 General Sentiment Analysis



Figure 5. General Sentiment Analysis of Trump-Biden (2020)

For this stage, the data was arranged in a pair of stacked bar charts, with each chart representing one figure or a pairing of figures. The data was arranged by sentiment as judged by VADER, with a tweet deemed as positive contributing to the green bar, a tweet deemed as neutral contributing to the orange bar, and negative to the red bar.

For each of these stacked bar charts, the x-axis arranged the tweets by week, with the date being the start of a given week. The left chart shows the total number of tweets that appeared each week, while the right chart shows the percentage of positive vs negative sentiment that appeared each week. This was necessary as the total number of tweets varied greatly week by week, and was hard to compare between bars using the left chart alone.

There were two elements to creating these charts. Firstly, iterating through each relevant data frame and creating a list for each week of how many positive, neutral and negative tweets there were as well as the total number. Secondly, taking those lists and using Matplotlib to create the stacked bar chart.

These charts were then automatically saved as PNGs for later referral, as generating the charts was rather time-consuming.

## 4.1.2 Hashtag Popularity

This stage was primarily used to help identify what was causing spikes in the data as displayed in the General Sentiment Analysis section.

I considered using word clouds for this section, as the requirement for this visualisation was to display words by their occurrence in the chosen week. However while this allowed me to quickly pick out the most popular hashtags, it was not suitable for exploring the prevalence of various conspiracy theories, which would be important later.

Instead, I settled for a simple table for each data frame, showing the top ten positive and negative hashtags for each week. I then applied a colour scheme to allow certain results to jump out easily.

The colour assignment was as follows; covid = red, qanon = green, Ukraine & Russia = blue.

1. Very relevant; a bright colour that indicates a strong correlation between the hashtag and a specific conspiracy theory.
2. Somewhat relevant; a pale colour that indicates a slight correlation between the hashtags and a specific conspiracy theory, but may be used in unrelated tweets.
3. No relevancy; any tweets that did not correlate with any of the three conspiracy theories remained white.
4. Generic conspiracy; any hashtag that indicated belief in a conspiracy theory not included in the selected three, was coloured white. Note that the list of 'generic' hashtags is not exhaustive, as the shared vernacular of conspiracy theorists is extensive.

This is achieved by creating two 2d arrays, one containing all of the necessary data that will populate the cells of the table, and one checking the data in each cell against a set of dictionaries to see if the cell should have a specific colour. A table is then generated from these two arrays using Matplotlib.

As the General Sentiment Analysis was only generating 8 graphs, each containing data from 15 data frames, it was fine to run them all at once, and easy to refer back to them individually. However, as this stage had a table for each figure/pairing and each data frame within that figure/pairing, making a chart for each data frame would have entailed the generation of 120 charts which was untenable.

Instead, I wrote three simple functions to allow the generation of a single chart by passing the desired figure/pairing and week, all charts for a specific figure/pairing, or all charts for a specific week.

## 4.1.3 Conspiracy Theory Prevalence
For this stage, the number of tweets that used at least one conspiracy-related hashtag was recorded for each conspiracy and was displayed similarly to the General Sentiment Analysis Charts, but with a grouped bar chart instead of a stacked bar chart.

## 4.1.4 Frequency Distribution
This step allowed for the checking of the frequency of individual words in a corpus made easy with the application of the NLTK FreqDist class.

To run a Frequency Distribution on the entire corpus would have been very time-consuming, so instead, I only worked with a single data frame at a time (that is a single week for a figure/pairing, i.e. 'Trump-Biden03-01').

Another similar process I explored and used to some extent was ngrams, specifically bigrams and trigrams. Instead of retrieving single words, this process returned a list of the most popular two or three long series of words.

While this was a very interesting tool, and caught the intentions of phrases used within data frames better than simple frequency distribution, I found it unnecessary for what I was trying to achieve.

## 4.1.5 Concordance

Concordance was a very useful tool for checking a specific word, often a hashtag in the case of this project and seeing the context is appears within its corpus. Again this process was made easy with the NLTK 'Text' object and 'concordance' method.

## 4.2 Investigating the Data

The first objective was to examine the graphs that were generated from General Sentiment Analysis to try and notice any trends or peculiarities, and then explore if there was a preeminent reason for the changes.

## 4.2.1 Exploring the first Graphs

At this stage, the graphs were rather chaotic, with the distribution of sentiment varying so much from week to week not creating any sort of baseline to reliably compare outliers too. However, there were some clearly anomalous weeks.

As can be seen in Figure 6, there was a colossal positive spike in the week 05-10 for the Biden-Harris pairing.



Figure 6. General Sentiment Analysis of Biden-Harris (2020)

A similar event can be seen during the week 05-17 for Biden in Figure 7 below.

Figure 7. General Sentiment Analysis of Biden (2020)

An initial examination shows that with an increased number of tweets the sentiment distribution is more prone to spiking. The reason for this is clear, as having such a large increase of tweets of a certain sentiment will overpower the other sentiments. This is especially prevalent in data frames with a smaller total number of entries.

As can be seen in Figure 8 below, the over 6,000,000 tweets that constitute Trump's data frame caused a kind of smoothing of the sentiment distribution when compared to other graphs, as weeks where there were outliers of a thousand tweets in one direction or another would not cause such a strong effect.



Figure 8. General Sentiment Analysis of Trump (2020)

This was not the case across all the graphs, as while spikes in the total number of tweets often caused spikes in sentiment distribution, there were also situations where spikes in sentiment distribution were present without large spikes in the total number of tweets.

Such an event can be seen during the week 03-01 and 05-31 for Trump-Harris in Figure 9 below.

Figure 9. General Sentiment Analysis of Trump-Harris (2020)

While the higher number of total tweets explaining why the percentage changed, it was not clear why there was such a dramatic spike for those weeks specifically. Further exploration of that data frame was required.

## 4.2.2 Investigating the data frame

I tried several different methods for investigating the tweets that appeared in the Biden-Harris data frame.

## 4.2.2.1 Popular Hashtags

Using the Hashtag Popularity Table for this data frame showed that the huge increase in sentiment matched with a hashtag in wide use for that data frame; '#obamagate'. This was not the case when attempting this approach across other data frames, as sometimes spikes in sentiment were not related to any particular hashtag.

Biden-Harris's Most Popular Hashtags in 05-10

| Positive | Neutral | Negative |
|---|---|---|
| ('#obamagate', 2647) | ('#bidenharris2020', 2) | ('#obamagate', 4) |
| N/A | N/A | ('#demrat', 2) |
| N/A | N/A | ('#bidenharris2020', 2) |
| N/A | N/A | ('#tarareade', 1) |
| N/A | N/A | ('#maga', 1) |
| N/A | N/A | ('#qanon', 1) |
| N/A | N/A | ('#bidenisobamagatefallguy', 1) |
| N/A | N/A | ('#neverpotus', 1) |
| N/A | N/A | ('#resistance', 1) |
| N/A | N/A | N/A |

Figure 10. Popular Hashtags for Biden-Harris 05-10 (2020)

At this point, some issues with the VADER sentiment analyser can start to be seen. It seems that the '#obamagate' hashtag is causing the positive spike in this data

frame, but it makes little sense for a conspiracy theory about Obama spying on Trump to be positive for Biden or Harris.

Similarly, there is a large positive spike in the 'Biden05-17' data frame; even though this is the week he claimed to any African-American voter that didn't vote for him that "you ain't black", which caused a considerable amount of controversy. Even though those hashtags can be seen in the table, they did not seem to affect the overall sentiment distribution for that week.

Biden's Most Popular Hashtags in 05-17

| Positive | Neutral | Negative |
|---|---|---|
| ('#obamagate', 41) | ('#obamagate', 124) | ('#obamagate', 350) |
| ('#biden', 33) | ('#biden', 98) | ('#youaintblack', 347) |
| ('#lindseygraham', 25) | ('#biden2020', 95) | ('#america', 290) |
| ('#burr', 25) | ('#deletefacebook', 92) | ('#votered', 290) |
| ('#gowdy', 25) | ('#youaintblack', 92) | ('#bidenkneweverything', 103) |
| ('#mcconnell', 25) | ('#joebidenisaracist', 79) | ('#obamakneweverything', 103) |
| ('#democrats', 18) | ('#joebiden', 53) | ('#lockthemallupnow', 103) |
| ('#biden2020', 17) | ('#ukraine', 49) | ('#factsmatter', 91) |
| ('#hillary', 16) | ('#pharma', 45) | ('#nogoquidprojoe', 90) |
| ('#nancypelosi', 16) | ('#swinefluvaccinefraud1976', 45) | ('#votealldemsout', 90) |

**Figure 11. Popular Hashtags for Biden 05-17 (2020)**

## 4.2.2.2 Frequency Distribution

For data frames where there was no obvious cause from the Hashtag Popularity Table, further investigation could be performed using frequency distribution.

Creating two corpora, one for the Biden-Harris 05-10 and one for the Biden 05-17 data frames with stopwords removed, then running them through the FreqDist 'most_common' method shows the following top ten words along with their number of occurrences.

| Biden-Harris 05-10 | Freq. | Biden 05-17 | Freq. |
|---|---|---|---|
| obama | 2650 | biden | 57390 |
| kamala | 2650 | joe | 56822 |
| perfect | 2649 | campaign | 56388 |
| shadow | 2648 | disinformation | 56250 |
| deep | 2648 | win | 56246 |
| vp | 2648 | china | 56234 |
| state | 2647 | done | 56214 |
| subversive | 2647 | continue | 56193 |
| maneuvering | 2647 | race | 56191 |
| framing | 2647 | massive | 56189 |

**Table 2. Frequency Distribution of two data frames**

From the number of words with near identical occurrences, it seems quite clear that these spikes are not caused by organic discussion on Twitter, but rather the result of a tweet that received a huge amount of retweets.

## 4.2.2.3 Concordance

Creating two further corpora, similar to the previous two but with stopwords remaining, and using words gathered in the previous stages allows us to see the specific retweet that has caused the spike for each group. The first tweet was retweeted an inordinate number of times, no doubt because it was a tweet from

```
"China is on a massive disinformation campaign because they are
desperate to have Sleepy Joe Biden win the presidential race so
they can continue to rip-off the United States, as they have done
  for decades, until I came along!" 56178 retweets
```

```
"Obama with his Deep State subversive maneuvering and framing of
innocents wants Biden's VP to be Kamala Harris. She's the perfect
 kind of corrupt prosecutor to add to his list of speed-dialing
 shadow gov cronies. The Constitution has its framers. Obama has
 his too.
 #ObamaGate" 2647 retweets
```

Donald Trump.

This shows two important things to consider when exploring the data; the effect of retweets and the overall accuracy of the VADER sentiment analysis.

## 4.3.1 Retweets

As it was clear that retweets had a major effect on the General Sentiment Analysis, I added the ability to generate the graphs with and without retweets. As can be seen by comparing figures 12 and 13 below the change was quite substantial.

As could be expected, removing the retweets removed the majority of the spikes in sentiment distribution across all of the graphs. This effect was more pronounced in the data frames with fewer overall entries, as a tweet that is retweeted thousands of times has a far greater effect, as can be seen on the Harris 05-10 data frame (3172 entries) versus the Trump 05-10 data frame (654,651 entries).

The reduction in size amongst the figures varied from the total number of tweets being 83% of the original size down to 28%. The two largest drops were in the Trump and Biden data frames, both of which dropped down to 28%. For these two major figures, the majority of the discussion around them was not original discussion, but rather retweeting. A similar pattern was observed with an analysis of political discourse for the 2016 US Election when they reported that almost 70% of tweets in their datasets were retweets (Yaqub et al 2017).

**Figure 12. General Sentiment Analysis of Trump-Biden, excluding retweets (2020)**



**Figure 13. General Sentiment Analysis of Trump-Biden, including retweets (2020)**

What is more surprising is that removing retweets flattened the distribution of sentiment to be more or less uniform across each figure/pairing. One would expect that across a period of three and a half months, there would be shifts in the sentiment distribution depending on what was occurring, especially during such a vitriolic period as the 2020 election.

The topic of whether including or excluding the retweets gives a better understanding of the data itself and has pros and cons for both approaches.

Other studies have removed retweets as part of the initial cleaning (Mirani and Sasi 2016). Others have done so as a step to remove clusters of tweets that were one tweet that had been consistently retweeted (Godfrey et al 2014) or to prevent inflating the number of tweets (Ordun et al 2020).

However, the reasoning behind removing retweets is usually to reduce the corpus to include only tweets of 'original authorship', which is not suitable for a project tackling conspiracy theories.

The ability to retweet and the ease with which a specific message can be amplified from a single user to a wide audience is part of how conspiracy theories can spread so easily. One of the highlighted tweets above is from Donald Trump who received 56,178 retweets.

While those may not have been original thoughts from those people who retweeted, they still decided to spread that message, which contributed to the overall sentiment that existed on Twitter.

## 4.3.2 VADER Sentiment Analysis

After seeing the two tweets in the Concordance section above as being interpreted as having positive sentiment, further examination of how the VADER sentiment analysis tool worked was required.

VADER uses a lexicon and rules-based approach to sentiment analysis (Hutto & Gilbert 2014), where the input text is compared against a large lexicon of 7520 words/emojis which assigns an intensity value of -4 for strongly negative to 4 for strongly positive to each of the words in the lexicon.

Other factors are also taken into account when calculating the sentiment score for a piece of text, such as capitilisation, punctuation (exclamation marks and question marks), words that negate the proceeding word, and words that boost the proceeding word.

With this in mind and looking at the tweet to see which words from the VADER lexicon were picked up, we can see some issues.

| Tweet text | VADER lexicon words (score) |
|---|---|
| Obama with his Deep State [subversive] maneuvering and framing of [innocents] wants Biden's VP to be Kamala Harris. She's the [perfect] [kind] of corrupt prosecutor to add to his list of speed-dialing shadow gov cronies. The Constitution has its framers. Obama has his too.  #ObamaGate | Subversive (-0.9) |
| | Innocents (1.1) |
| | Perfect (2.7) |
| | Kind (2.4) |
| Tweet text | VADER lexicon words (score) |
| China is on a massive disinformation campaign because they are [desperate] to have Sleepy Joe Biden [win] the presidential race so they can continue to rip-off the [United] States, as they have done for decades, until I came along[!] | desperate (-1.3) |
| | win (2.8) |
| | United (1.8) |
| | ! (Applies a 0.292 intensity amplifier to sentiment) |

Table 3. Breakdown of VADER sentiment analysis

In the case of the first tweet from Table 3, there is some unfortunate use of the words 'perfect kind', which is used negatively in the tweet, but as they are both very positive words it throws off VADER's sentiment analysis.

This example also illustrates issues that make it hard for NLP to work in general, but is only exacerbated by the microblogging format of Twitter. Due to the more casual nature of communication by tweet, misspellings, contractions, and peculiar slang is common. In the example above 'maneuvering' is a misspelling, and 'gov' is a contraction of government.

The second tweet on the table above includes an exclamation mark that amplifies any sentiment which is present. It also shows a quirk that any mention of the United States will add to the tweet's positive sentiment, which has potentially far-reaching consequences in a study concerning a United States election.

While neither of these would have made a huge difference to the final sentiment analysis, they show some of the problems that can potentially exist across all tweets (Giachanou & Crestani 2016).

It seems that out of the box, while still very powerful, VADER is not perfectly equipped to interpret tweets related to conspiracy theories. While it picks up on "subversive" in the first tweet, there is a large amount of conspiracy theory vernacular that would undoubtedly help VADER to correctly interpret this tweet.  The addition of words such as 'framing, 'corrupt' & 'cronies' would help to cut through the confusion that 'perfect kind' has caused in this tweet.

Additionally, VADER has implementation for sentiment-laden idioms, which could be used for popular conspiracy phrases such as 'deep state' & 'shadow gov'.

The VADER documentation on GitHub (Hutto 2021) describes the rigorous process that was used to determine the value assigned to each word. To follow the same process it is recommended to "find 10 independent humans to evaluate/rate each new token you want to add to the lexicon, make sure the standard deviation doesn't exceed 2.5 and take the average rating for the valence. This will keep the file consistent."

This process would have far exceeded the scope of the project. A less rigorous approach could have been pursued, by assigning an arbitrary number for each word that 'felt right. But beyond being rather unempirical, the sheer extent and complexity of a conspiracy theorist's vocabulary made even attempting the task very time-consuming.

Therefore I left the VADER lexicon as it was, despite some of its flaws for the task at hand.

I experimented with an alternative way of presenting the data but did not pursue it in great detail due to a lack of time. Instead of splitting the sentiment analysis into positive, neutral & negative, I would only use positive & negative categories.

As the nature of the corpus is specifically calling out misinformation, it makes sense that most of the discourse would be very polarising with a more negative slant. It is hard to imagine a tweet discussing how the deep state is trying to undermine Trump's election campaign to be neutral in its sentiment.

### 4.3.3 Conspiracy Theory Prevalence

The final means by which I explored the data was checking just how prevalent each of the three conspiracy theories was within each figure/pairing overall. Perhaps unsurprisingly, tweets that mentioned Trump alone were most reliably full of conspiracy theory chatter.

As can be seen on the left-hand side of the tables below, Trump has by far the most conspiracy theory discussion in the tweets which are associated with him alone. He has the highest in Covid & qanon conspiracy theories, while Biden has the highest in Ukraine conspiracy theories, which makes sense as he is a central figure in that belief.

Furthermore, when comparing the data frames of Biden, Pence, and Harris against the data frames where they are mentioned alongside Trump (Trump-Biden, Trump-Pence, Trump-Harris), the percentage of tweets discussing conspiracy theories generally increases by up to 153% in most cases, with one anomalous result (Trump-Pence, QAnon; 3.59%) nearly being an order of magnitude larger than its solo result (Pence, QAnon; 0.34%).

| Conspiracy Prevalence Including Retweets | | | | | Covid | QAnon | Ukraine |
|---|---|---|---|---|---|---|---|
| | Covid | QAnon | Ukraine | | Covid | QAnon | Ukraine |
| Trump | 1.14% | 1.90% | 0.44% | N/A | - | - | - |
| Biden | 0.18% | 0.69% | 0.80% | Trump-Biden | 0.14% | 1.11% | 2.02% |
| Pence | 0.28% | 0.34% | 0.07% | Trump-Pence | 0.63% | 3.59% | 0.14% |
| Harris | 0.29% | 0.39% | 0.13% | Trump-Harris | 0.24% | 0.99% | 0.32% |

Table 4. Conspiracy Prevalence, including retweets (2020)

This large spike is due to a large retweet spike in that data frame. The table below shows the same layout of data as above, but excluding retweets.

The trend of the Trump tweets containing the most conspiracy theory related hashtags remains when excluding retweets. As does the increase of conspiracy theory related hashtags as Trump is added to pairings of other figures, with the exception of the solo Biden data frames having more on the Ukraine conspiracy theories.

Another thing to note is that when excluding retweets, the percentage of tweets containing conspiracy theory related hashtags for Trump goes down, whereas for Biden, Pence, and Harris this percentage generally increases. This implies that many retweets in the Trump data frame are conspiracy related.

For the other three, the increased percentage implies that many of the retweets in those data frames are not conspiracy related, and the use of conspiracy related hashtags originates from 'original authorship'.

This shows part of the dangerous power that social media provides, one that was wielded by Trump and his supporters to great effect, using retweets to amplify a message.

| Conspiracy Prevalence Excluding Retweets | | | | | | | |
|---|---|---|---|---|---|---|---|
| | Covid | QAnon | Ukraine | | Covid | QAnon | Ukraine |
| Trump | 0.65% | 1.32% | 0.31% | N/A | - | - | - |
| Biden | 0.20% | 0.85% | 0.72% | Trump-Biden | 0.31% | 0.88% | 0.52% |
| Pence | 0.48% | 0.61% | 0.07% | Trump-Pence | 0.75% | 0.76% | 0.25% |
| Harris | 0.46% | 0.52% | 0.22% | Trump-Harris | 0.46% | 1.14% | 0.57% |

Table 5. Conspiracy Prevalence, excluding retweets (2020)

## 4.4 Overall Sentiment Analysis

When looking at each figure/pairing as a whole, and not divided by week we can see the overall sentiment distribution for each figure/pairing. On the whole, there are several similarities with the previous two tables.

Namely, tweets discussing only Trump have the highest percentage of negative sentiment, and when paired up with the other figures the negative percentage tends to increase compared to the figure's solo sentiment distribution.

| Sentiment Distribution Including Retweets | | | | | | | |
|---|---|---|---|---|---|---|---|
| | Pos. | Neu. | Neg. | | Pos. | Neu. | Neg. |
| Trump | 8.76% | 43.43% | 47.80% | N/A | - | - | - |
| Biden | 27.25% | 36.97% | 35.77% | Trump-Biden | 10.72% | 45.54% | 43.71% |
| Pence | 7.57% | 55.39% | 36.96% | Trump-Pence | 11.65% | 45.53% | 42.76% |
| Harris | 6.76% | 58.73% | 34.45% | Trump-Harris | 5.85% | 58.27% | 24.89% |

Table 6. Sentiment Distribution, including retweets (2020)

| Sentiment Distribution Excluding Retweets | | | | | | | |
|---|---|---|---|---|---|---|---|
| | Pos. | Neu. | Neg. | | Pos. | Neu. | Neg. |
| Trump | 9.01% | 42.30% | 48.68% | N/A | - | - | - |
| Biden | 11.69% | 46.82% | 41.47% | Trump-Biden | 12.18% | 40.17% | 47.64% |
| Pence | 10.46% | 48.11% | 41.42% | Trump-Pence | 10.10% | 38.54% | 51.35% |
| Harris | 10.11% | 48.33% | 41.55% | Trump-Harris | 10.10% | 39.87% | 50.01% |

Table 7. Sentiment Distribution, excluding retweets (2020)

## 5.1 2022 Data

On August 8[th] 2022 Donald Trump's Mar-a-Lago residence was raided by the FBI, as part of a  federal criminal investigation into the former president. This unsurprisingly caused an uproar from Trump's militant supporters, with some going so far as to picket FBI offices and in one case attempt an attack (NPR 2022).

As it is want to do, the landscape of the social media world has changed quite a bit in the past two years, and such a major event in the story of the central figure of the qanon conspiracy theory poses a perfect opportunity to examine how it may have changed.

One of the primary ways in which the qanon conspiracy theory presence is likely to have changed between 2020 and 2022 is that after the January 6[th] Insurrection, many Twitter accounts associated with qanon were banned from the platform.

Those users migrated to other platforms such as Parler and Truth Social, so while there is undoubtedly still a qanon presence on Twitter, it will most likely be severely reduced.

## 5.1.1 Checking Conspiracy Prevalence

With the data from 2022, I started by looking at the Conspiracy Prevalence. This was because I expected that over the course of 2 years, the vernacular used by these conspiracy communities would have evolved and the dictionary of hashtags I had generated for the 2020 data may not be as effective.

Below is the 2022 data, I only included the data with retweets included for this section along with the percentage increases and decreases when compared to 2020.

|  | Corpus size | Covid | QAnon | Ukraine |
|---|---|---|---|---|
| Trump | 749,399 | 0.06% (-1.08%) | 0.60% (-1.30%) | 0.38% (-0.06%) |
| Biden | 259,078 | 0.15% (-0.03%) | 0.91% (+0.21%) | 0.77% (+0.08%) |
| Pence | 14,582 | 0.09% (-0.19%) | 0.52% (+0.18%) | 0.05% (-0.02%) |
| Harris | 11,916 | 0.01% (-0.28%) | 0.20% (+0.19%) | 0.03% (-0.04%) |
| Trump-Biden | 56,615 | 0.05% (-0.09%) | 0.65% (-0.46%) | 1.65% (-0.37%) |
| Trump-Pence | 10,513 | 0.01% (-0.62%) | 0.16% (-3.43%) | 0.08% (-0.06%) |
| Trump-Harris | 555 | 0% (-0.24%) | 1.08% (+0.09%) | 0.72% (+0.40%) |

Table 8. Conspiracy Prevalence, including retweets (2022)

The prevalence of covid conspiracy tweets across all data frames has greatly decreased. Where the lowest prevalence in 2020 was 0.14% in the Trump-Biden data frames, the highest prevalence in 2022 is 0.15% in the Biden data frames.

This is very likely because during 2020 the pandemic was at its height, while now the impact on people's lives has decreased and so there is less pushback exhibiting as conspiracy theorists.

Additionally, with the exception of the popular '#hunterbidenslaptop' hashtag that was picked up in the Biden07-06 data frame, the Ukraine conspiracy theory has taken a similar decrease in popularity. As this particular conspiracy was likely

Results from the Trump-Harris data frame are rather random, but that can be explained by the fact that it is far too small to draw any reasonable conclusions, coming in at only 555 tweets.

In fact, all data frames besides Trump, Biden, and Trump-Biden were rather thin, all coming in below 15,000 entries. So I decided for the 2022 data to focus on the Trump, Biden, and Trump-Biden data frames, and only on the QAnon conspiracy within those data frames.

## 5.1.2 Finding New Hashtags

So I set out to see if there was a set of new hashtags that the QAnon conspiracy community was using on Twitter, using the same methodology I did when first generated my lists of conspiracy theory related hashtags.

I attempted to locate a list of any current popular QAnon hashtags but was unable to find anything, so decided to work from my pre-existing list. While not ideal, some of the old hashtags would still be seeing some use and would hopefully provide links to new hashtags.

However, running the old list of hashtags through the Hashtag Extractor proved far less useful than expected.

```
Tags from Initial List:  [('#qanon', 199), ('#pizzagate', 72), ('#maga', 54), ('#trumpwon', 23), ('#lgb', 21), ('#letsgobrando
n', 21), ('#greatawakening', 21), ('#gqp', 15), ('#wwg1wga', 11), ('#qanons', 11), ('#linwood', 5), ('#thegreatawakening', 2),
('#thestorm', 2), ('#redpill78', 1), ('#wearethenewsnow', 1), ('#pedogate', 1)]

Extracted Tags:  [('#demvoice1', 362), ('#bluevoices', 362), ('#trublue', 362), ('#trump', 122), ('#trumpcult', 52), ('#illumin
ati', 39), ('#qanoncult', 35), ('#chemtrails', 34), ('#gop', 32), ('#soros', 29), ('#qdupes', 29), ('#birthers', 28), ('#truthe
rs', 28), ('#plandemic', 27), ('#biden', 23), ('#factsmatter', 21), ('#antivaxxers', 21), ('#newworldorder', 19), ('#sethrich',
19), ('#benghazi', 19), ('#obamagate', 19), ('#newsupdate', 18), ('#donaldtrump', 18), ('#sharpiegate', 16), ('#smartnews', 1
5), ('#truth', 15), ('#chavez', 15), ('#hunterbiden', 14), ('#democrats', 14), ('#conspiracytheory', 14), ('#lizardpeople', 1
4), ('#gopcoverup', 13), ('#freedom', 13), ('#conspiracytheories', 13), ('#uranium1', 13), ('#thebiglie', 12), ('#4chan', 12),
('#nwo', 12), ('#covid', 12), ('#january6thhearings', 12), ('#cult45', 12), ('#mindcontrol', 12), ('#jfk', 12), ('#45s', 12),
('#flatearth', 11), ('#wakeup', 11), ('#truthseeker', 11), ('#votebluetosavedemocracy', 11), ('#antivaxxer', 11), ('#vaccinechi
ps', 11), ('#goodbyegop', 10), ('#roevwade', 10), ('#scotus', 10), ('#usa', 10), ('#trumpcultists', 10), ('#fascism', 9), ('#bi
glie', 9), ('#epstein', 9), ('#ghislainemaxwellclientlist', 9), ('#wherearethechildren', 9)]
```

**Figure 14. Initial List of Hashtags (top) Extract Hashtags from 2022 Data (bottom)**

There is a very small number of hashtags from the existing conspiracy hashtag list that have appeared across the 2022 data. Comparing this to results from the 2020 data and accounting for the difference in size, this constitutes a %%% decrease in QAnon hashtag use.

What's more, a large number of these extracted hashtags have no relation to QAnon. The tags '#demvoice1', '#bluevoices', and '#trueblue', are from one tweet that received several retweets that merely mentioned QAnon. Then a great deal of what may be expected to be QAnon or at least conspiracy related such as '#truthers', '#sharpiegate, '#soros', '#deepstate', '#lizardpeople', is again from a

single retweet calling trump cultists deluded, followed by a mocking string of 18 ridiculous conspiracy tags.

Some of the remaining tags such as '#gqp' (a portmanteau of GOP and Q, referring to Republican politicians supporting QAnon conspiracy theories) and '#trumpwon' imply QAnon membership. I also opened up to more generic hashtags that may not have QAnon implications, such as '#maga','#trumpwon','#letsgobrandon', and '#lgb'.

Despite adding these to the list of hashtags and running it through the Hashtag Extractor again, it returned no results that could be linked to QAnon communities.

It is quite likely that there are bubbles of QAnon communities that were either hiding in the data or that were simply not picked up in the data gathering. However, possibly a better explanation is that QAnon conversation has simply left Twitter in favour of other social media platforms such as Parler and Truth Social, which are more accepting of those kinds of beliefs.

## 6.1 Conclusions

The goals of this project as set were rather broad, and while there is much refinement that could be done and further research to be completed, I believe they were all met.

Firstly, through the use of various python libraries such as os and gzip I was able to handle the data within the huge misinformation callout corpus, and organise it into many reasonably sized data frames with pandas. From these data frames, I could easily select the appropriate data frame/s for the problem at hand with the powerful tools afforded me by pandas.

Within the data itself, there were several things found that distinguished the presence and prevalence of conspiracy theories and specific sentiments between the various figures and pairings. Perhaps the largest, and the most unsurprising, was that of the Trump data frame. It had the highest percentage of Negative sentiment, as well as the highest percentage of conspiracy related tweets. This carried over to any data frame which was another figure who was paired with Trump, as his inclusion in a pairing caused both negative sentiment and conspiracy prevalence to rise in comparison to the original figure's solo data frame.

I became acquainted with the VADER sentiment analysis tool and found it to be a very useful tool. Part of what makes it so useful is not just the gold-standard lexicon that the creators have laboriously created, but also the ease with which somebody could make additions to the lexicon to make it fit their purpose.

## 6.2 Future Work

For future work, there is a wide breadth of directions I could envision taking this project just within the limits of the approaches taken over the course of this project.

It would be interesting to expand the scope of the data explored. An earlier starting point would cover the Democratic Primaries and could benefit from looking at other prominent figures, the controversial figure of Bernie Sanders would undoubtedly reveal some things of interest.

A later endpoint could stretch up to the election, which could be very interesting to see the discourse grow as Trump's defeat loomed and Biden was finally victorious. Qanon activity would probably grow, and while the dataset used in this project did not have much to go on regarding accusations of election fraud, that kind of discourse was very common leading up to election day and the immediate aftermath.

An even later date, which could be the most interesting of all would have been to look at the January 6th Insurrection, as this was the point where followers of the qanon conspiracy theory truly came into the limelight. The immediate aftermath would also be interesting to see, as there was a major culling of qanon Twitter

accounts after the fact (BBC News 2021), so it could be researched how effective that was in curtailing qanon discourse on the platform.

Another time period to be explored would be the 2016 election, and see how the prevalence of conspiracy theories between the two differed or were similar. The pizzagate conspiracy that existed during the 2016 election was a precursor to qanon after all.

Other social media platforms could be interesting to explore. While places like Reddit and 4chan (or other imageboards) have proven to be havens for conspiracy theorists, platforms such as Parler and Truth Social would be very interesting to look at.  These are platforms that the alt-right have adopted after their exile from Twitter, with the latter being founded by Trump himself.

In terms of technical improvements that could be made that I was unable to pursue in the scope of this project, further exploration of sentiment analysis would be interesting. One avenue to pursue would be adding to the VADER lexicon to improve the analysis of conspiracy theories, or potentially exploring other forms of sentiment analysis.

Another interesting area that I would have liked to try, but did not have the proper time to cover in the appropriate depth, is that of machine learning. It can be applied to NLP, with research done showing ML models being very effective at identifying conspiracy theories (Marcinello et al 2021).


## 6.3 Reflection

The process of developing this project has been challenging, made none the easier that I chose to engage in an area of study with which I had very little experience previously.

Throughout the course of this project, I have learned a great deal having greatly developed my Python skills as well as coming to grips with the exciting field of Natural Language Processing and dipping my toes into Machine Learning.

I've also become well acquainted with the pandas library, and am confident in creating and manipulating data frames. I also had to learn to manage and explore huge quantities of data.

However, I feel the organisation or structure of the development of my project could have benefitted from more forethought and planning. While the topic of the 2020 US election was of interest to me, I felt rather unsure of the direction to take the overall project.

I think if I had established a more concrete objective earlier in the research process the overall result of my project could have greatly benefitted. I felt drawn in many directions with my research, and as such do not feel I took a deep enough study of

any particular area. Still, failures in the organisation and planning a project taught me important lessons in self-management in projects of this nature for the future.

**References**

1. BBC News. 2022. *Twitter suspends 70,000 accounts linked to QAnon*. [online] Available at: https://www.bbc.co.uk/news/technology-55638558 [Accessed 22 September 2022].

2. Bird, S., Klein, E. and Loper, E., 2009. *Natural language processing with Python*. Beijing: O'Reilly. Liddy, E.D. 2001. Natural Language Processing. In Encyclopedia of Library and Information Science, 2nd Ed. NY. Marcel Decker, Inc.

3. Chowdhary, K.R. (2020). Natural Language Processing. In: Fundamentals of Artificial Intelligence. Springer, New Delhi. https://doi.org/10.1007/978-81-322-3972-7_19

4. Dilley, L., Welna, W. & Foster, F. (2021/2022). QAnon Propaganda on Twitter as Information Warfare: Influencers, Networks, and Narratives. Accepted Sept. 16, 2021 at Frontiers in Communication, 6:707595, doi: 10.3389/fcomm.2021.707595

5. Facebook. 2022. *2022 Q2 Report*. [online] Available at: https://s21.q4cdn.com/399680738/files/doc_financials/2022/q2/Meta-06.30.2022-Exhibit-99.1-Final.pdf [Accessed 20 September 2022].

6. Giachanou, A. and Crestani, F., 2016. Like It or Not. *ACM Computing Surveys*, 49(2), pp.1-41.

7. Godfrey et al. (2014). A Case Study in Text Mining: Interpreting Twitter Data From World Cup Tweets.

8. Hunter, J., 2007. Matplotlib: A 2D Graphics Environment. *Computing in Science & Engineering*, 9(3), pp.90-95.

9. Hutto, C., n.d. *GitHub VADER Sentiment Analysis.* [online] GitHub. Available at: https://github.com/cjhutto/vaderSentiment [Accessed 22 September 2022].

10. Hutto, C. and Gilbert, E. (2014) VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text. *Proceedings of the International AAAI Conference on Web and Social Media*, 8(1), pp. 216-225. Available at: https://ojs.aaai.org/index.php/ICWSM/article/view/14550 [Accessed: 22 September 2022]

11. Institute for Strategic Dialogue. 2020. *The Genesis of a Conspiracy Theory.* Available at: https://www.isdglobal.org/wp-content/uploads/2020/07/The-Genesis-of-a-Conspiracy-Theory.pdf

12. Khan, G., Swar, B. and Lee, S., 2014. Social Media Risks and Benefits. *Social Science Computer Review*, 32(5), pp.606-627.

13. Mahase, E., 2020. Hydroxychloroquine for covid-19: the end of the line?. *BMJ*, p.m2378.

14. Marcellino, William, Todd C. Helmus, Joshua Kerrigan, Hilary Reininger, Rouslan I. Karimov, and Rebecca Ann Lawrenc. 2020, *Detecting Conspiracy Theories on Social Media: Improving Machine Learning to Detect and Understand*. Online Conspiracy Theories, Santa Monica, Calif.: RAND Corporation, RR-A676-1, 2021. [Accessed 22, 2022] Available at: https://www.rand.org/pubs/research_reports/RRA676-1.html

15. Marwick, A. and Lewis, R. 2017. *Media Manipulation and Disinformation Online.* Data & Society, pp35-36.

16. Mirani, T.B., & Sasi, S. (2016). *Sentiment Analysis of ISIS Related Tweets Using Absolute Location.* 2016 International Conference on Computational Science and Computational Intelligence (CSCI), 1140-1145.

17. NPR.org. 2022. *An attempted attack on an FBI office raises concerns about violent far-right rhetoric*. [online] Available at: https://www.npr.org/2022/08/12/1117275044/an-attempted-attack-on-an-fbi-office-raises-concerns-about-violent-far-right-rhe [Accessed 22 September 2022].

18. Ordun, Catherine & Purushotham, Sanjay & Raff, Edward. (2020). Exploratory Analysis of Covid-19 Tweets using Topic Modeling, UMAP, and DiGraphs. Available at: https://arxiv.org/pdf/2005.03082.pdf Accessed 22 September 2022]

19. The pandas development team, 2020. *pandas-dev/pandas: Pandas*.

20. Preece, A., Spasic, I., Evans, K., Rogers, D., Webberley, W., Roberts, C. and Innes, M., 2018. Sentinel: A Codesigned Platform for Semantic Enrichment of Social Media Streams. *IEEE Transactions on Computational Social Systems*, 5(1), pp.118-131.

21. Roesslein, J., 2020. Tweepy: Twitter for Python! URL: *https://github.com/tweepy/tweepy.*

22. Shao, C., Ciampaglia, G., Varol, O., Yang, K., Flammini, A. and Menczer, F., 2018. The spread of low-credibility content by social bots. *Nature Communications*, 9(1).

23. Stecula, D. and Pickup, M., 2021. Social Media, Cognitive Reflection, and Conspiracy Beliefs. *Frontiers in Political Science*, 3.

24. Twitter, 2022. *Twitter 2022 Q2 Report*. [online] S22.q4cdn.com. Available at: https://s22.q4cdn.com/826641620/files/doc_financials/2022/q2/Final_Q2'22_Earnings_Release.pdf [Accessed 22 September 2022].

25. Wardle, C. and Derakhshan, H. 2017. *Information disorder: Toward an interdisciplinary framework for research and policymaking.* Pp. 21-22, 49-56. Available at: http://tverezo.info/wp-content/uploads/2017/11/PREMS-162317-GBR-2018-Report-desinformation-A4-BAT.pdf [Accessed 22 September 2022].

26. Yaqub, U., Chun, S., Atluri, V. and Vaidya, J., 2017. Analysis of political discourse on twitter in the context of the 2016 US presidential elections. *Government Information Quarterly*, 34(4), pp.613-626.

# Sentiment Analysis from Conspiracy Theories in the 2020 US Election

*Robert Bryant*

**MSc Computing**

**Professor Alun Preece**

**School of Computer Science and Informatics, Cardiff University**

**23/09/2022**

**Abstract**

In the 2020 US Election there was a huge amount of discussion on Twitter, much of it hostile and much of it based on conspiracy. I will use Natural Language Processing, and specifically Sentiment Analysis to interrogate and analyse a large corpus of Twitter data from the period.

I first had to distil the large Twitter dataset into parts that were easier to process, split by prominent figures (Trump, Biden, Pence, Harris). These data frames were then run through various pipelines to explore the Sentiment Analysis and Conspiracy Prevalence.

The analysis showed that Trump as a figure trended towards negative sentiment and conspiracy prevalence, going so far as to increase the negative sentiment and conspiracy prevalence of those he was paired with. Additionally, retweets were found to have a major effect on sentiment and conspiracies.

There were various questions answered, primarily concerning the negative influence of Trump. Some areas could be improved, specifically the accuracy of the Sentiment Analysis tools and a more accurate means to determine tweets with specific conspiracy theories.

# Contents

## Table of Figures

## Table of Tables

## 1.1 Introduction

Social media as technology has grown since its inception in the 90s with the widespread adoption of the World Wide Web. From the beginnings of relatively simple sites such as SixDegrees.com and GeoCities, monoliths of the industry have emerged such as Facebook and Twitter. From each company's Second Quarter 2022 Results, Twitter reports 237.8 million daily users (Twitter 2022) while Facebook reports a staggering 1.97 billion daily users (Facebook 2022).

While these platforms have great benefits such as social connectivity, social involvement, information attainment, and entertainment, there are also many risks associated with their use (Khan et al 2014). These risks can include overuse, mental health issues, social problems, and privacy (ibid.). This study will focus on the societal drawbacks that can be enabled through the use of social media, using a set of Twitter data gathered from March to July of 2020, exploring the US election of that year and the prevalence of conspiracy theories in that period.

In the 2020 election, there was a huge amount of disinformation being propagated on social media, this was in no small part due to the controversial figure of the incumbent American President Trump who promoted much of the disinformation and conspiracy theories.

Using powerful tools included with the Natural Language ToolKit (NLTK) for Natural Language Processing (NLP), I hope to explore the prevalence of specific conspiracy theories during the election in relation to specific candidates or pairings of candidates, as well as any trends in sentiment analysis with these candidates or pairings of candidates.

## 1.2 Aims and Objectives

- Develop a tool for curating a large corpus of tweets and turning it into a manageable data set that fits my requirements.
- Explore which figures and pairings of figures (Trump, Biden, Pence, Harris) are most commonly associated with conspiracy theories/certain sentiment analyses.
- Evaluation of the VADER sentiment analysis tool for interpreting conspiracy theories on Twitter.

# 2 Background
## 2.1 Previous Work
### 2.1.1 Conspiracy Theories

There are three main conspiracy theories explored in this project; covid-related, QAnon, and Ukraine & Russia.

1. Conspiracy theories related to covid-19 are very broad and come in several varieties. There are those who believe covid is a bioweapon engineered in China, those who think the elite are using the opportunity to perform unnecessary vaccinations for some nefarious purpose, those who think the whole pandemic is a lie and that the disease doesn't actually exist or is seriously exaggerated. While not strictly a conspiracy theory while President, Trump contributed to spreading a great deal of misinformation including the use of dubious medical treatments (Mahase 2020. Additionally, the 2020 election was taking place at the height of the first wave of the pandemic.

   The wanton spread of this conspiracy theory throughout social media and society at large resulted in damage to public health, as certain groups of people would not engage in social distancing or other preventative measures and later refused to receive the vaccine (Romer & Jamieson 2020).

2. QAnon is a conspiracy theory that emerged on a message board and describes a wide-reaching conspiracy being perpetrated by the elite, in which Donald Trump is a sort of saviour. The theory is too complicated and bizarre to explain succinctly but with its close ties to Trump it has an important place in this project. Outside the scope of the data used in this project, QAnon had a huge impact in the real world as the so-called 'qarmy' were heavily involved in the January 6[th] Insurrection in 2021.

   Perhaps more than any other conspiracy theory QAnon has benefitted from the spread of social media. Not only was it born on the image board 4chan, but from its inception to June 17 2020 the Institute for Strategic Dialogue recorded "69,475,451 million tweets, 487,310 Facebook posts, and 281,554 Instagram posts mentioning QAnon-related

   hashtags and phrases" (ISD 2020).

3. The third conspiracy theory I have chosen to examine is regarding Joe Biden and the allegations that he had used his position as Vice-President to pressure Ukraine into dropping a corruption investigation that would have put his son Hunter Biden under scrutiny. This was pushed as an attempt to discredit Biden's presidential bid.
While not

## 2.1.2 Conspiracy Theories on Social Media

Conspiracy theories as a whole are a complex thing to define, with the task not growing any easier when trying to pin down a specific theory. One common thread of most conspiracy theories is the belief in a truth that does not align with the mainstream and often factual truth, often with the belief that some other party is suppressing the theorist's 'truth' for some nefarious reason (Marcellino et al 2021).

This leads to an entrenching of their position as the believers reject information that does not align with their theories and seek out sources that reinforce their beliefs. This process has been made much more prevalent with the use of social media. As a user of social media can curate who they follow, echo chambers can emerge. An echo chamber is a community where dissenting voices are not heard and the members' beliefs can be reinforced (Wardle & Derakhshan 2017).

Social media is a medium for which a great deal of misinformation can be spread, with people who use social media as a source of news increasing the chance that they will hold beliefs in a conspiracy theory (Stecula & Pickup 2021).

There are many ways that conspiracy theories can be spread across social media:

- Bots can be used by those within the conspiracy community or outside agents who wish to fan the flames of that theory (Shao et al 2018)
- Organised raids can be used to artificially inflate the reach of rumours or misinformation with reposting, resharing, or retweeting as well as general public discussion (Wardle & Derakhshan 2017)
- Hashtags can be forced to trend by concerted efforts by a large community, or an existing hashtag can be hijacked by posting irrelevant or harmful information using an already popular hashtag, with the intention of diluting the

content of the hashtag or with malicious intent for those participating in the conversation (Marwick & Lewis 2017)

### 2.1.3 Natural Language Processing

Natural Language Processing (NLP) is a wide-ranging computational technique and research domain, where the goal is for the comprehension of natural human language by a machine along with performing some level of linguistic analysis (Liddy 2010). This can be attempted and achieved with varying levels of success through a range of techniques.

A machine capable of understanding the idiosyncrasies of human speech is a tall task, as there is a wide range of quirks that humans simply take for granted, such as synonyms, idioms, and humour.

## 2.2 Python Libraries

As this project was written entirely in Python using Jupyter Notebooks, all of the modules and libraries utilised are Python modules and libraries.

### 2.2.1 File Management

os - This module allows for the use of operating system functionality. This project makes use of the listdir function to assist in opening all the files in a given directory.

gzip - This module allows for the compression and decompression of gzip files. This project makes use of the open function to read each individual tweet stored within the many gzip files provided.

### 2.2.4 Natural Language Processing

NLTK - The Natural Language Tool Kit is a suite of libraries for NLP. It is suitable for those who are first experimenting with NLP as well as practitioners and those conducting research into NLP (Chowdhary 2020), so it is no doubt ideal for a first project in the field.

An incredibly important suite, it includes all the base tools to perform NLP tasks that are useful to this project such as tokenization, removal of stop words, and with the

VADER model a great sentiment analysis tool that sits at the heart of much of this project. (Bird, Klein & Loper 2009)

VADER - Vader is a sentiment analysis tool that takes a lexicon and rules-based approach, specifically designed for use with social media posts. Additionally, VADER uses intensity scores (ranking just how positive or negative something is) rather than just a binary (deciding if something is either good or bad), which allows for a more nuanced result for sentiment analysis (Hutto & Gilbert 2014).

As sentiment analysis is at the very heart of this project, it would have been hard to get anything completed without an NLP tool of some kind and VADER was the one that fit the job best. Using the 'polarity_scores' method on a string returns a dictionary containing four scores; positive, neutral, negative & compound (a combination of the first three).

Though it stretched beyond the scope of this project, the tool could have been adapted to perform even better by adding additional conspiracy-related words with matching sentiment scores to the VADER lexicon.

## 2.2.3 Data Manipulation & Visualisation

pandas - pandas is a very powerful library allowing for data manipulation and analysis. One of the most essential libraries for this project, as it was used to create all of the data frames for the project (pandas 2020, McKinney 2010).

Using simple and powerful tools within pandas, the data could then be manipulated and used in concert with Matplotlib to generate data visualisations.

Additionally using the 'to_pickle' method, the many large data frames that were pulled from the misinformation callout corpus were able to be saved, so the extraction only needed to be performed once.

Matplotlib - Matplotlib is a library for the creation of a wide range of data visualisations. While the charts created in this project were all relatively rudimentary, the library is a powerful tool with the potential for creating animated and interactive visualisations (Hunter 2007).

As the project involved analysing millions of tweets, it would have been a great hindrance to do so without a good means of displaying the data in a more visual medium.

### 2.2.7 Tweepy

Tweepy is an easy-to-use library to access the Twitter API and gather the tweets that Twitter has made publicly available (Roesslein 2020).

As I was only engaging in a small amount of gathering of my own Twitter data, Tweepy was a perfect fit as I did not have to learn to use the Twitter API for a relatively small part of the overall project.

# 3 Data Collection and Organisation

## 3.1 Existing Corpus

The majority of data that is used in this study was taken from a huge dataset of Tweets ordered by month and year, ranging from 2019 to 2022. Within each month was a json.gz file, a compressed json file that contained millions of tweets. The data was gathered searching for tweets for which the text included any of a list of search terms relating to disinformation, misinformation, and conspiracy theories amongst other things. This was done using the Sentinel platform (Preece et al 2017).

Due to this expansive source of Twitter data, the necessity for my own data collection was quite minimal. However, for Part 5 there was some use of the Twitter API to collect tweets from August-September 2022, as the misinformation callout corpus did not include tweets from that period.

The misinformation callout corpus comprised of tweets gathered every day in the months of March, April, May & June 2020. The June set was missing entries from $15^{th}$-$28^{th}$, so I decided to simply drop the last few days of June as well, meaning I used data stretching over a total of 15 weeks.

All the json files in this period contained over 8 trillion tweets, so a careful selection would have to be made. This was not only due to the prohibitively large amount of data, but also the fact that the vast majority of it would have been irrelevant to the subject at hand.

## 3.2 Curating the Corpus

### 3.2.1 Dividing the Corpus by Important Figures

I decided to focus on important figures in the 2020 election, selecting both presidential candidates, Trump and Biden, and their running mates, Pence and Harris. From the main datasets tweets would be selected that mentioned either one or two of these figures, and then organised into weekly datasets for easier management.

The option of examining tweets that mentioned either three of the figures or all four together was explored, but the sizes of those curated datasets were all too small to be worth exploring further. Additionally, when I first gathered my own data frames

from the original corpus, I organised the data by month but was later changed to be organised by week to be able to more carefully examine the data.

As can be seen in Figure 1 and Figure 2 when organised by month, outlying data points would obfuscate a quarter of the data. While organised by week it gave a more precise impression of what was occurring, with the outlier only affecting a single week.
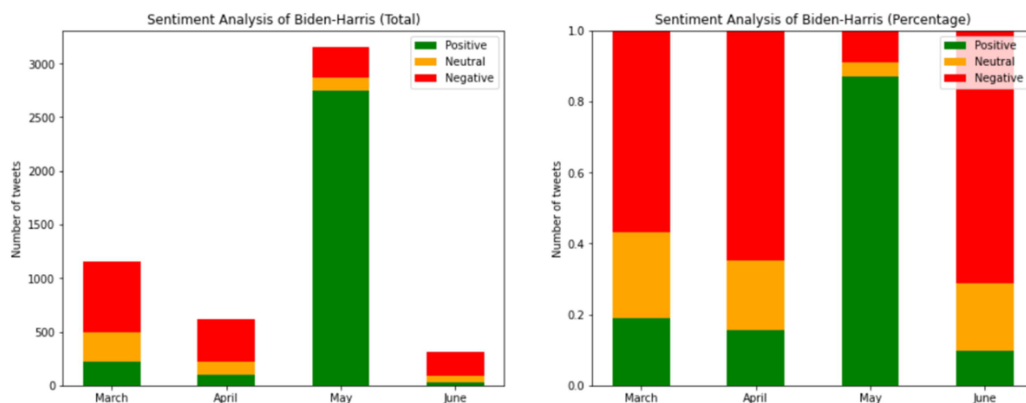


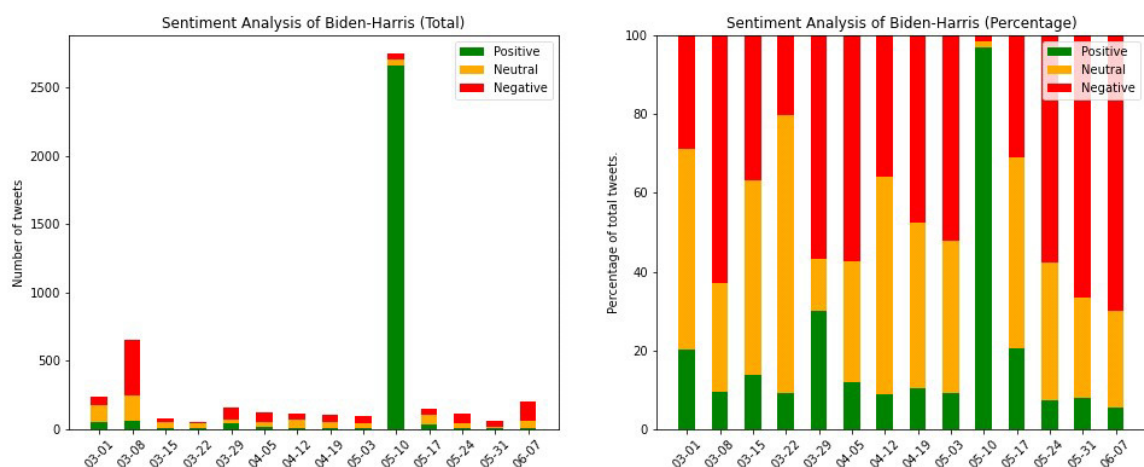Figure 1. General Sentiment Analysis of Biden-Harris (by month)



Figure 2. General Sentiment Analysis of Biden-Harris (by week)

Once each of the json.gz was unzipped each of these files ranged in sizes from 2-5 gigabytes, which was far too large to allow unzipping the entire dataset and working from that. Even unzipping the files by month proved to be a cumbersome and untenable process.

Instead, I iterated through the directories the files were saved to and used the gzip library 'open' method to read each file within each directory without saving it. While this was still time-consuming, it saved a huge amount of disk space.

### 3.2.2 Columns for the Data Frame

The corpus had a huge range of potential data points for every tweet contained within, so for my own datasets I chose to focus on the following fields, some of which were taken directly from the corpus and some of which were derived from data within it;

1. user_id: the id of the individual user who made the tweet. Note that the intention of this was to check for 'super users' whose tweets are prevalent in the dataset, and not as an attempt to identify the anonymous users.

2. tweet_id: the id of the individual tweet.

3. text: the content of the tweet. From this field, the names, hashtags, sentiment, conspiracies, trumpTags & bidenTags were derived.
   Certain common pre-processing techniques such as removing punctuation and lower casing were not performed, as these elements are important parts of the VADER sentiment analysis process.
   Mentions (any token beginning with '@') and URLs (any token beginning with 'http') were removed, as they contributed

4. date: the date and time the tweet was made

5. names: which of the important figures were mentioned in this tweet

6. hashtags: any hashtags in the tweet, ignoring individual '#'s

7. sentimentA: the compound sentiment score as provided by VADER, converted to either 'Positive' (if score is >= 0.5), 'Neutral' (if score is > -0.5 and < 0.5) or 'Negative' (if score is <=-0.5

8. sentimentB: an alternate form of recording the sentiment. Still using VADER, but instead of calculating based on the compound score, simply recording whether the positive score or negative score was higher.

9. retweet_count: the number of retweets the tweet received, if any

10. conspiracies: if any hashtags that were associated with a specific conspiracy theory were mentioned in the body of the tweet, then the name of that conspiracy theory was added to this list

11. trumpTags: an integer showing how many Trump hashtags were mentioned in the tweet
12. bidenTags: an integer showing how many Biden hashtags were mentioned in the tweet

Deciding on these columns was an iterative process, as the focus and scope of the project expanded more columns were added. For instance, initially the data that appears in the final 'conspiracies' column was being pulled from hashtags while performing NLP in later sections. Moving the generation of the hashtags and making it a permanent column of the data frame made it easier to explore the data later on.

### 3.2.3 Unused Fields
There were some fields included in the data frames which I had intended to utilise, but in the end did not have enough time to implement them in a meaningful way.

- User ID was going to be used to find if there were any 'super users' who posted more regularly and in doing so skewed the data, and see if 'super users' were more common in a specific figure/pairing or conspiracy theory. However more things proved more informative and this angle was dropped.
- TrumpTags and BidenTags were the first steps in trying to organise tweets by Trump supporters and Biden supporters, in a similar way to how conspiracy theories were organised by hashtag. Unfortunately, this approach was far too rudimentary, and to properly realise this idea would have required a heavy investment of time and machine learning.
- SentimentB was intended to be used with the General Sentiment Analysis to show sentiment distributions of Positive/Negative instead of Positive/Neutral/Negative. While I did create these graphs, I did not find that they added a significant amount of further analysis so did not include them in this report.

## 3.3 Categorising Hashtags
### 3.3.1 Hashtag Generation
To create the 'conspiracies' field for my dataset I needed some way to identify which tweets were talking about specific conspiracy theories.

As I was working with such a large dataset, I chose a relatively simple solution. As I already planned to gather the hashtags that were present in each tweet I decided to utilise that data, as hashtags are good indicators of a tweet's overall message.

n.b. I also utilised a similar workflow to populate the bidenTags and trumpTags fields, but as they ended up being unused in the final project I won't describe it in detail.

### 3.3.2 Choosing Conspiracy Theories

I chose some theories which I thought may have been prevalent during this period, with an initial selection of; covid, QAnon, Ukraine (the Hunter-Biden scandal), Russia (Russian interference in the election), fraud (voter fraud), and climate change denial.

Through performing the following steps a few things became clear that led to a refinement of these topics.

1. Many of the hashtags used in the Ukraine and Russia conspiracy hashtags were identical, so I combined them into one topic.
2. While there were large accusations of voter fraud during the 2020 election that did not begin until the time surrounding Election Day itself, and that was over four months after the end of the available data. Due to this, the number of relevant tweets was rather low and I removed them.
3. While climate change denial is endemic on Twitter, it was not a major topic of discussion about the 2020 election when compared to other conspiracies and so I also removed climate change denial from the running.

This left covid conspiracies, which were rampant as it was still the height of the pandemic; QAnon conspiracies, which were similarly popular due to their intrinsic link to Donald Trump; and the renamed ukr_ru category which was popular due to the linked event's impacts on the 2020 election.

### 3.3.3 False Positives

Hashtags are by no means a perfect method of categorising tweets. While some hashtags are used by conspiracy theorists, they are just as likely to not be. Something like '#maga' fits in this category, as while it is commonly used within the QAnon conspiracy theory, it is also used widely by Trump supporters as a whole,

and not all tweets supporting Trump are also discussing QAnon. Additionally, somebody could be using a hashtag that they do not agree with if the tweet's intention is derision or if they are attempting to 'hijack' a hashtag.

Due to this, I was careful to only select hashtags which had a strong bias in them implicitly. For instance '#chinavirus', '#ccpvirus', and '#chinaliedpeopledied' all indicate a belief that China was behind the covid-19 pandemic. Additionally, there were a great number of strange, almost coded hashtags which were in wide use by members of the QAnon conspiracy theory such as '#thestorm','#wwg1wga', and'#thegreatawakening'.

Targeting such pointed or coded tweets such as these reduced the likelihood that somebody was merely discussing the conspiracy theory. These conclusions were drawn by using concordance to check the use of hashtags in context within their tweets.

### 3.3.4 Hashtag Extraction Methodology

To generate the final lists of conspiracy hashtags I created initial lists from two sources;

1. Looking at the most popular overall hashtags from 'curated dataset a' and taking any hashtags that were clearly related to a specific conspiracy theory, backed up by exploring them with concordance.
2. Finding lists of hashtags that had been compiled in similar research which listed hashtags commonly used by the specific conspiracy theories I followed (Dilley, Welna & Foster 2021).

With these initial lists, I then ran the entire curated dataset through a workflow that created a dictionary for each conspiracy theory recording the number of occurrences of any hashtags that appeared alongside my initial list.

This new list required a bit of curating as they often included hashtags that were either not intrinsically linked to the conspiracy theory or were not relevant at all. To accomplish this curation I took hashtags from the new list that I wanted to check and looked up any occurrences of the hashtag in the corpus using concordance from NLTK.

I then combined the initial list with the new list and ran the combined list through the whole process again. I only had to perform this process twice before no new hashtags were found.



Figure 3. Hashtag Extraction Workflow

## 3.4 Corpus Results

Once I ran the large dataset through this workflow, I was left with a total of 7,445,214 tweets, distributed between the individual and paired figures as follows;

| Trump | Trump-Biden | Trump-Pence | Trump-Harris | Biden |
|---|---|---|---|---|
| 6,441,479 | 276,013 | 32,789 | 9,953 | 491,877 |
| Biden-Harris | Biden-Pence | Pence | Pence-Harris | Harris |
| 5,121 | 503 | 26,882 | 32 | 34,709 |

Table 1. Size of figure/pairing data frames

Due to the low number of tweets in the Biden-Pence & Pence-Harris datasets, I chose to remove them, as they would be unlikely to show any interesting results.

## 3.5 Gathering 2022 Data

On August 8[th] 2022 Trump's Mar-a-Lago estate was raided by the FBI in search of classified documents that had allegedly been illegally retained after he failed to be re-elected for a second term. This caused an outcry from Trump supporters, so it seemed a relevant topic to explore and compare the state of tweets currently to the tweets during the 2020 elections.

In the misinformation callout dataset there had been tweets gathered for June & July 2022, but nothing beyond those months, so I chose to gather my own data using the Twitter API and Tweepy library.

## 3.5.1 Twitter API

The Twitter API allows access to all extant public tweets and a great deal of data associated with each Tweet. I applied for the Twitter Developer Account and was granted the necessary API keys to begin making requests to the API and gathering data. As I only had the Essential Access Level available to me I was limited to a maximum of 500,000 tweets per month and only being able to access tweets from the previous 7 days.

Due to this the data I personally gathered was not as extensive as the misinformation callout dataset. I managed to gather tweets from August 5[th]-September 1st with 25,000 tweets per day.

## 3.5.2 Tweepy

I used the Tweepy library for easy access to the Twitter API. To try and keep my data as similar to the other dataset as possible I used as many of the same search terms as I could. Due to the Essential Access Level, there were limits on the length of the search query string and I had to truncate the list. The only missing terms however were related to Europe, e.g. 'european False', 'european Fraud', 'european Hoax'.

As I would be further curating these tweets to split the gathered corpus into the same important figures that the 2020 dataset is split into, the removal of these European tags seemed inconsequential as the chances of them being related to any of the US figures were low.

Due to user error and misunderstanding of the various fields available to be pulled from the Twitter API, this set of data does not record how many retweets a tweet received, or if it was a retweet at all.

### 3.5.3 Most Popular Times

While it did not make much of a difference given the relatively small amount of data I was able to gather, I set out to determine the best time of day to gather tweets. This was simply done by iterating through all the data frames and incrementing the correct hour in an array.

The resulting graphs showed the most popular times as between 10 pm-3 am UTC, which would be 5 pm-10 pm EST or 2 pm-7 pm PST.

**Figure 4. Popular Tweeting Times**

## 4 2020 Election Dataframes

This data was explored in many ways and went through many changes throughout the project, but can be broadly described in the following ways;

1. General Sentiment Analysis
2. Hashtag Popularity
3. Conspiracy Theory Prevalence
4. Frequency Distribution
5. Concordance

## 4.1.1 General Sentiment Analysis



Figure 5. General Sentiment Analysis of Trump-Biden (2020)

For this stage, the data was arranged in a pair of stacked bar charts, with each chart representing one figure or a pairing of figures. The data was arranged by sentiment as judged by VADER, with a tweet deemed as positive contributing to the green bar, a tweet deemed as neutral contributing to the orange bar, and negative to the red bar.

For each of these stacked bar charts, the x-axis arranged the tweets by week, with the date being the start of a given week. The left chart shows the total number of tweets that appeared each week, while the right chart shows the percentage of positive vs negative sentiment that appeared each week. This was necessary as the total number of tweets varied greatly week by week, and was hard to compare between bars using the left chart alone.

There were two elements to creating these charts. Firstly, iterating through each relevant data frame and creating a list for each week of how many positive, neutral and negative tweets there were as well as the total number. Secondly, taking those lists and using Matplotlib to create the stacked bar chart.

These charts were then automatically saved as PNGs for later referral, as generating the charts was rather time-consuming.

## 4.1.2 Hashtag Popularity

This stage was primarily used to help identify what was causing spikes in the data as displayed in the General Sentiment Analysis section.

I considered using word clouds for this section, as the requirement for this visualisation was to display words by their occurrence in the chosen week. However while this allowed me to quickly pick out the most popular hashtags, it was not suitable for exploring the prevalence of various conspiracy theories, which would be important later.

Instead, I settled for a simple table for each data frame, showing the top ten positive and negative hashtags for each week. I then applied a colour scheme to allow certain results to jump out easily.

The colour assignment was as follows; covid = red, qanon = green, Ukraine & Russia = blue.

1. Very relevant; a bright colour that indicates a strong correlation between the hashtag and a specific conspiracy theory.
2. Somewhat relevant; a pale colour that indicates a slight correlation between the hashtags and a specific conspiracy theory, but may be used in unrelated tweets.
3. No relevancy; any tweets that did not correlate with any of the three conspiracy theories remained white.
4. Generic conspiracy; any hashtag that indicated belief in a conspiracy theory not included in the selected three, was coloured white. Note that the list of 'generic' hashtags is not exhaustive, as the shared vernacular of conspiracy theorists is extensive.

This is achieved by creating two 2d arrays, one containing all of the necessary data that will populate the cells of the table, and one checking the data in each cell against a set of dictionaries to see if the cell should have a specific colour. A table is then generated from these two arrays using Matplotlib.

As the General Sentiment Analysis was only generating 8 graphs, each containing data from 15 data frames, it was fine to run them all at once, and easy to refer back to them individually. However, as this stage had a table for each figure/pairing and each data frame within that figure/pairing, making a chart for each data frame would have entailed the generation of 120 charts which was untenable.

Instead, I wrote three simple functions to allow the generation of a single chart by passing the desired figure/pairing and week, all charts for a specific figure/pairing, or all charts for a specific week.

## 4.1.3 Conspiracy Theory Prevalence

For this stage, the number of tweets that used at least one conspiracy-related hashtag was recorded for each conspiracy and was displayed similarly to the General Sentiment Analysis Charts, but with a grouped bar chart instead of a stacked bar chart.

## 4.1.4 Frequency Distribution

This step allowed for the checking of the frequency of individual words in a corpus made easy with the application of the NLTK FreqDist class.

To run a Frequency Distribution on the entire corpus would have been very time-consuming, so instead, I only worked with a single data frame at a time (that is a single week for a figure/pairing, i.e. 'Trump-Biden03-01').

Another similar process I explored and used to some extent was ngrams, specifically bigrams and trigrams. Instead of retrieving single words, this process returned a list of the most popular two or three long series of words.

While this was a very interesting tool, and caught the intentions of phrases used within data frames better than simple frequency distribution, I found it unnecessary for what I was trying to achieve.

## 4.1.5 Concordance

Concordance was a very useful tool for checking a specific word, often a hashtag in the case of this project and seeing the context is appears within its corpus. Again this process was made easy with the NLTK 'Text' object and 'concordance' method.

## 4.2 Investigating the Data

The first objective was to examine the graphs that were generated from General Sentiment Analysis to try and notice any trends or peculiarities, and then explore if there was a preeminent reason for the changes.

## 4.2.1 Exploring the first Graphs

At this stage, the graphs were rather chaotic, with the distribution of sentiment varying so much from week to week not creating any sort of baseline to reliably compare outliers too. However, there were some clearly anomalous weeks.

As can be seen in Figure 6, there was a colossal positive spike in the week 05-10 for the Biden-Harris pairing.



Figure 6. General Sentiment Analysis of Biden-Harris (2020)

A similar event can be seen during the week 05-17 for Biden in Figure 7 below.

Figure 7. General Sentiment Analysis of Biden (2020)

An initial examination shows that with an increased number of tweets the sentiment distribution is more prone to spiking. The reason for this is clear, as having such a large increase of tweets of a certain sentiment will overpower the other sentiments. This is especially prevalent in data frames with a smaller total number of entries.

As can be seen in Figure 8 below, the over 6,000,000 tweets that constitute Trump's data frame caused a kind of smoothing of the sentiment distribution when compared to other graphs, as weeks where there were outliers of a thousand tweets in one direction or another would not cause such a strong effect.



Figure 8. General Sentiment Analysis of Trump (2020)

This was not the case across all the graphs, as while spikes in the total number of tweets often caused spikes in sentiment distribution, there were also situations where spikes in sentiment distribution were present without large spikes in the total number of tweets.

Such an event can be seen during the week 03-01 and 05-31 for Trump-Harris in Figure 9 below.

Figure 9. General Sentiment Analysis of Trump-Harris (2020)

While the higher number of total tweets explaining why the percentage changed, it was not clear why there was such a dramatic spike for those weeks specifically. Further exploration of that data frame was required.

## 4.2.2 Investigating the data frame

I tried several different methods for investigating the tweets that appeared in the Biden-Harris data frame.

## 4.2.2.1 Popular Hashtags

Using the Hashtag Popularity Table for this data frame showed that the huge increase in sentiment matched with a hashtag in wide use for that data frame; '#obamagate'. This was not the case when attempting this approach across other data frames, as sometimes spikes in sentiment were not related to any particular hashtag.

### Biden-Harris's Most Popular Hashtags in 05-10

| Positive | Neutral | Negative |
|---|---|---|
| ('#obamagate', 2647) | ('#bidenharris2020', 2) | ('#obamagate', 4) |
| N/A | N/A | ('#demrat', 2) |
| N/A | N/A | ('#bidenharris2020', 2) |
| N/A | N/A | ('#tarareade', 1) |
| N/A | N/A | ('#maga', 1) |
| N/A | N/A | ('#qanon', 1) |
| N/A | N/A | ('#bidenisobamagatefallguy', 1) |
| N/A | N/A | ('#neverpotus', 1) |
| N/A | N/A | ('#resistance', 1) |
| N/A | N/A | N/A |

Figure 10. Popular Hashtags for Biden-Harris 05-10 (2020)

At this point, some issues with the VADER sentiment analyser can start to be seen. It seems that the '#obamagate' hashtag is causing the positive spike in this data

frame, but it makes little sense for a conspiracy theory about Obama spying on Trump to be positive for Biden or Harris.

Similarly, there is a large positive spike in the 'Biden05-17' data frame; even though this is the week he claimed to any African-American voter that didn't vote for him that "you ain't black", which caused a considerable amount of controversy. Even though those hashtags can be seen in the table, they did not seem to affect the overall sentiment distribution for that week.

Biden's Most Popular Hashtags in 05-17

| Positive | Neutral | Negative |
|---|---|---|
| ('#obamagate', 41) | ('#obamagate', 124) | ('#obamagate', 350) |
| ('#biden', 33) | ('#biden', 98) | ('#youaintblack', 347) |
| ('#lindseygraham', 25) | ('#biden2020', 95) | ('#america', 290) |
| ('#burr', 25) | ('#deletefacebook', 92) | ('#votered', 290) |
| ('#gowdy', 25) | ('#youaintblack', 92) | ('#bidenkneweverything', 103) |
| ('#mcconnell', 25) | ('#joebidenisaracist', 79) | ('#obamakneweverything', 103) |
| ('#democrats', 18) | ('#joebiden', 53) | ('#lockthemallupnow', 103) |
| ('#biden2020', 17) | ('#ukraine', 49) | ('#factsmatter', 91) |
| ('#hillary', 16) | ('#pharma', 45) | ('#nogoquidprojoe', 90) |
| ('#nancypelosi', 16) | ('#swinefluvaccinefraud1976', 45) | ('#votealldemsout', 90) |

Figure 11. Popular Hashtags for Biden 05-17 (2020)

## 4.2.2.2 Frequency Distribution

For data frames where there was no obvious cause from the Hashtag Popularity Table, further investigation could be performed using frequency distribution.

Creating two corpora, one for the Biden-Harris 05-10 and one for the Biden 05-17 data frames with stopwords removed, then running them through the FreqDist 'most_common' method shows the following top ten words along with their number of occurrences.

| Biden-Harris 05-10 | Freq. | Biden 05-17 | Freq. |
|---|---|---|---|
| obama | 2650 | biden | 57390 |
| kamala | 2650 | joe | 56822 |
| perfect | 2649 | campaign | 56388 |
| shadow | 2648 | disinformation | 56250 |
| deep | 2648 | win | 56246 |
| vp | 2648 | china | 56234 |
| state | 2647 | done | 56214 |
| subversive | 2647 | continue | 56193 |
| maneuvering | 2647 | race | 56191 |
| framing | 2647 | massive | 56189 |

Table 2. Frequency Distribution of two data frames

From the number of words with near identical occurrences, it seems quite clear that these spikes are not caused by organic discussion on Twitter, but rather the result of a tweet that received a huge amount of retweets.

## 4.2.2.3 Concordance

Creating two further corpora, similar to the previous two but with stopwords remaining, and using words gathered in the previous stages allows us to see the specific retweet that has caused the spike for each group. The first tweet was retweeted an inordinate number of times, no doubt because it was a tweet from

```
"China is on a massive disinformation campaign because they are
desperate to have Sleepy Joe Biden win the presidential race so
they can continue to rip-off the United States, as they have done
  for decades, until I came along!" 56178 retweets
```

```
"Obama with his Deep State subversive maneuvering and framing of
innocents wants Biden's VP to be Kamala Harris. She's the perfect
 kind of corrupt prosecutor to add to his list of speed-dialing
 shadow gov cronies. The Constitution has its framers. Obama has
 his too.
 #ObamaGate" 2647 retweets
```

Donald Trump.

This shows two important things to consider when exploring the data; the effect of retweets and the overall accuracy of the VADER sentiment analysis.

## 4.3.1 Retweets

As it was clear that retweets had a major effect on the General Sentiment Analysis, I added the ability to generate the graphs with and without retweets. As can be seen by comparing figures 12 and 13 below the change was quite substantial.

As could be expected, removing the retweets removed the majority of the spikes in sentiment distribution across all of the graphs. This effect was more pronounced in the data frames with fewer overall entries, as a tweet that is retweeted thousands of times has a far greater effect, as can be seen on the Harris 05-10 data frame (3172 entries) versus the Trump 05-10 data frame (654,651 entries).

The reduction in size amongst the figures varied from the total number of tweets being 83% of the original size down to 28%. The two largest drops were in the Trump and Biden data frames, both of which dropped down to 28%. For these two major figures, the majority of the discussion around them was not original discussion, but rather retweeting. A similar pattern was observed with an analysis of political discourse for the 2016 US Election when they reported that almost 70% of tweets in their datasets were retweets (Yaqub et al 2017).

**Figure 12. General Sentiment Analysis of Trump-Biden, excluding retweets (2020)**



**Figure 13. General Sentiment Analysis of Trump-Biden, including retweets (2020)**

What is more surprising is that removing retweets flattened the distribution of sentiment to be more or less uniform across each figure/pairing. One would expect that across a period of three and a half months, there would be shifts in the sentiment distribution depending on what was occurring, especially during such a vitriolic period as the 2020 election.

The topic of whether including or excluding the retweets gives a better understanding of the data itself and has pros and cons for both approaches.

Other studies have removed retweets as part of the initial cleaning (Mirani and Sasi 2016). Others have done so as a step to remove clusters of tweets that were one tweet that had been consistently retweeted (Godfrey et al 2014) or to prevent inflating the number of tweets (Ordun et al 2020).

However, the reasoning behind removing retweets is usually to reduce the corpus to include only tweets of 'original authorship', which is not suitable for a project tackling conspiracy theories.

The ability to retweet and the ease with which a specific message can be amplified from a single user to a wide audience is part of how conspiracy theories can spread so easily. One of the highlighted tweets above is from Donald Trump who received 56,178 retweets.

While those may not have been original thoughts from those people who retweeted, they still decided to spread that message, which contributed to the overall sentiment that existed on Twitter.

## 4.3.2 VADER Sentiment Analysis

After seeing the two tweets in the Concordance section above as being interpreted as having positive sentiment, further examination of how the VADER sentiment analysis tool worked was required.

VADER uses a lexicon and rules-based approach to sentiment analysis (Hutto & Gilbert 2014), where the input text is compared against a large lexicon of 7520 words/emojis which assigns an intensity value of -4 for strongly negative to 4 for strongly positive to each of the words in the lexicon.

Other factors are also taken into account when calculating the sentiment score for a piece of text, such as capitilisation, punctuation (exclamation marks and question marks), words that negate the proceeding word, and words that boost the proceeding word.

With this in mind and looking at the tweet to see which words from the VADER lexicon were picked up, we can see some issues.

| Tweet text | VADER lexicon words (score) |
|---|---|
| Obama with his Deep State [subversive] maneuvering and framing of [innocents] wants Biden's VP to be Kamala Harris. She's the [perfect] [kind] of corrupt prosecutor to add to his list of speed-dialing shadow gov cronies. The Constitution has its framers. Obama has his too. #ObamaGate | Subversive (-0.9) |
| | Innocents (1.1) |
| | Perfect (2.7) |
| | Kind (2.4) |
| Tweet text | VADER lexicon words (score) |
| China is on a massive disinformation campaign because they are [desperate] to have Sleepy Joe Biden [win] the presidential race so they can continue to rip-off the [United] States, as they have done for decades, until I came along[!] | desperate (-1.3) |
| | win (2.8) |
| | United (1.8) |
| | ! (Applies a 0.292 intensity amplifier to sentiment) |

Table 3. Breakdown of VADER sentiment analysis

In the case of the first tweet from Table 3, there is some unfortunate use of the words 'perfect kind', which is used negatively in the tweet, but as they are both very positive words it throws off VADER's sentiment analysis.

This example also illustrates issues that make it hard for NLP to work in general, but is only exacerbated by the microblogging format of Twitter. Due to the more casual nature of communication by tweet, misspellings, contractions, and peculiar slang is common. In the example above 'maneuvering' is a misspelling, and 'gov' is a contraction of government.

The second tweet on the table above includes an exclamation mark that amplifies any sentiment which is present. It also shows a quirk that any mention of the United States will add to the tweet's positive sentiment, which has potentially far-reaching consequences in a study concerning a United States election.

While neither of these would have made a huge difference to the final sentiment analysis, they show some of the problems that can potentially exist across all tweets (Giachanou & Crestani 2016).

It seems that out of the box, while still very powerful, VADER is not perfectly equipped to interpret tweets related to conspiracy theories. While it picks up on "subversive" in the first tweet, there is a large amount of conspiracy theory vernacular that would undoubtedly help VADER to correctly interpret this tweet. The addition of words such as 'framing, 'corrupt' & 'cronies' would help to cut through the confusion that 'perfect kind' has caused in this tweet.

Additionally, VADER has implementation for sentiment-laden idioms, which could be used for popular conspiracy phrases such as 'deep state' & 'shadow gov'.

The VADER documentation on GitHub (Hutto 2021) describes the rigorous process that was used to determine the value assigned to each word. To follow the same process it is recommended to "find 10 independent humans to evaluate/rate each new token you want to add to the lexicon, make sure the standard deviation doesn't exceed 2.5 and take the average rating for the valence. This will keep the file consistent."

This process would have far exceeded the scope of the project. A less rigorous approach could have been pursued, by assigning an arbitrary number for each word that 'felt right. But beyond being rather unempirical, the sheer extent and complexity of a conspiracy theorist's vocabulary made even attempting the task very time-consuming.

Therefore I left the VADER lexicon as it was, despite some of its flaws for the task at hand.

I experimented with an alternative way of presenting the data but did not pursue it in great detail due to a lack of time. Instead of splitting the sentiment analysis into positive, neutral & negative, I would only use positive & negative categories.

As the nature of the corpus is specifically calling out misinformation, it makes sense that most of the discourse would be very polarising with a more negative slant. It is hard to imagine a tweet discussing how the deep state is trying to undermine Trump's election campaign to be neutral in its sentiment.

### 4.3.3 Conspiracy Theory Prevalence

The final means by which I explored the data was checking just how prevalent each of the three conspiracy theories was within each figure/pairing overall. Perhaps unsurprisingly, tweets that mentioned Trump alone were most reliably full of conspiracy theory chatter.

As can be seen on the left-hand side of the tables below, Trump has by far the most conspiracy theory discussion in the tweets which are associated with him alone. He has the highest in Covid & qanon conspiracy theories, while Biden has the highest in Ukraine conspiracy theories, which makes sense as he is a central figure in that belief.

Furthermore, when comparing the data frames of Biden, Pence, and Harris against the data frames where they are mentioned alongside Trump (Trump-Biden, Trump-Pence, Trump-Harris), the percentage of tweets discussing conspiracy theories generally increases by up to 153% in most cases, with one anomalous result (Trump-Pence, QAnon; 3.59%) nearly being an order of magnitude larger than its solo result (Pence, QAnon; 0.34%).

| Conspiracy Prevalence Including Retweets | | | | | Covid | QAnon | Ukraine |
|---|---|---|---|---|---|---|---|
| | Covid | QAnon | Ukraine | | Covid | QAnon | Ukraine |
| Trump | 1.14% | 1.90% | 0.44% | N/A | - | - | - |
| Biden | 0.18% | 0.69% | 0.80% | Trump-Biden | 0.14% | 1.11% | 2.02% |
| Pence | 0.28% | 0.34% | 0.07% | Trump-Pence | 0.63% | 3.59% | 0.14% |
| Harris | 0.29% | 0.39% | 0.13% | Trump-Harris | 0.24% | 0.99% | 0.32% |

Table 4. Conspiracy Prevalence, including retweets (2020)

This large spike is due to a large retweet spike in that data frame. The table below shows the same layout of data as above, but excluding retweets.

The trend of the Trump tweets containing the most conspiracy theory related hashtags remains when excluding retweets. As does the increase of conspiracy theory related hashtags as Trump is added to pairings of other figures, with the exception of the solo Biden data frames having more on the Ukraine conspiracy theories.

Another thing to note is that when excluding retweets, the percentage of tweets containing conspiracy theory related hashtags for Trump goes down, whereas for Biden, Pence, and Harris this percentage generally increases. This implies that many retweets in the Trump data frame are conspiracy related.

For the other three, the increased percentage implies that many of the retweets in those data frames are not conspiracy related, and the use of conspiracy related hashtags originates from 'original authorship'.

This shows part of the dangerous power that social media provides, one that was wielded by Trump and his supporters to great effect, using retweets to amplify a message.

| Conspiracy Prevalence Excluding Retweets | | | | | | | |
|---|---|---|---|---|---|---|---|
| | Covid | QAnon | Ukraine | | Covid | QAnon | Ukraine |
| Trump | 0.65% | 1.32% | 0.31% | N/A | - | - | - |
| Biden | 0.20% | 0.85% | 0.72% | Trump-Biden | 0.31% | 0.88% | 0.52% |
| Pence | 0.48% | 0.61% | 0.07% | Trump-Pence | 0.75% | 0.76% | 0.25% |
| Harris | 0.46% | 0.52% | 0.22% | Trump-Harris | 0.46% | 1.14% | 0.57% |

Table 5. Conspiracy Prevalence, excluding retweets (2020)

## 4.4 Overall Sentiment Analysis

When looking at each figure/pairing as a whole, and not divided by week we can see the overall sentiment distribution for each figure/pairing. On the whole, there are several similarities with the previous two tables.

Namely, tweets discussing only Trump have the highest percentage of negative sentiment, and when paired up with the other figures the negative percentage tends to increase compared to the figure's solo sentiment distribution.

| Sentiment Distribution Including Retweets | | | | | | | |
|---|---|---|---|---|---|---|---|
| | Pos. | Neu. | Neg. | | Pos. | Neu. | Neg. |
| Trump | 8.76% | 43.43% | 47.80% | N/A | - | - | - |
| Biden | 27.25% | 36.97% | 35.77% | Trump-Biden | 10.72% | 45.54% | 43.71% |
| Pence | 7.57% | 55.39% | 36.96% | Trump-Pence | 11.65% | 45.53% | 42.76% |
| Harris | 6.76% | 58.73% | 34.45% | Trump-Harris | 5.85% | 58.27% | 24.89% |

Table 6. Sentiment Distribution, including retweets (2020)

| Sentiment Distribution Excluding Retweets | | | | | | | |
|---|---|---|---|---|---|---|---|
| | Pos. | Neu. | Neg. | | Pos. | Neu. | Neg. |
| Trump | 9.01% | 42.30% | 48.68% | N/A | - | - | - |
| Biden | 11.69% | 46.82% | 41.47% | Trump-Biden | 12.18% | 40.17% | 47.64% |
| Pence | 10.46% | 48.11% | 41.42% | Trump-Pence | 10.10% | 38.54% | 51.35% |
| Harris | 10.11% | 48.33% | 41.55% | Trump-Harris | 10.10% | 39.87% | 50.01% |

Table 7. Sentiment Distribution, excluding retweets (2020)

## 5.1 2022 Data

On August 8[th] 2022 Donald Trump's Mar-a-Lago residence was raided by the FBI, as part of a  federal criminal investigation into the former president. This unsurprisingly caused an uproar from Trump's militant supporters, with some going so far as to picket FBI offices and in one case attempt an attack (NPR 2022).

As it is want to do, the landscape of the social media world has changed quite a bit in the past two years, and such a major event in the story of the central figure of the qanon conspiracy theory poses a perfect opportunity to examine how it may have changed.

One of the primary ways in which the qanon conspiracy theory presence is likely to have changed between 2020 and 2022 is that after the January 6[th] Insurrection, many Twitter accounts associated with qanon were banned from the platform.

Those users migrated to other platforms such as Parler and Truth Social, so while there is undoubtedly still a qanon presence on Twitter, it will most likely be severely reduced.

## 5.1.1 Checking Conspiracy Prevalence

With the data from 2022, I started by looking at the Conspiracy Prevalence. This was because I expected that over the course of 2 years, the vernacular used by these conspiracy communities would have evolved and the dictionary of hashtags I had generated for the 2020 data may not be as effective.

Below is the 2022 data, I only included the data with retweets included for this section along with the percentage increases and decreases when compared to 2020.

|  | Corpus size | Covid | QAnon | Ukraine |
|---|---|---|---|---|
| Trump | 749,399 | 0.06% (-1.08%) | 0.60% (-1.30%) | 0.38% (-0.06%) |
| Biden | 259,078 | 0.15% (-0.03%) | 0.91% (+0.21%) | 0.77% (+0.08%) |
| Pence | 14,582 | 0.09% (-0.19%) | 0.52% (+0.18%) | 0.05% (-0.02%) |
| Harris | 11,916 | 0.01% (-0.28%) | 0.20% (+0.19%) | 0.03% (-0.04%) |
| Trump-Biden | 56,615 | 0.05% (-0.09%) | 0.65% (-0.46%) | 1.65% (-0.37%) |
| Trump-Pence | 10,513 | 0.01% (-0.62%) | 0.16% (-3.43%) | 0.08% (-0.06%) |
| Trump-Harris | 555 | 0% (-0.24%) | 1.08% (+0.09%) | 0.72% (+0.40%) |

Table 8. Conspiracy Prevalence, including retweets (2022)

The prevalence of covid conspiracy tweets across all data frames has greatly decreased. Where the lowest prevalence in 2020 was 0.14% in the Trump-Biden data frames, the highest prevalence in 2022 is 0.15% in the Biden data frames.

This is very likely because during 2020 the pandemic was at its height, while now the impact on people's lives has decreased and so there is less pushback exhibiting as conspiracy theorists.

Additionally, with the exception of the popular '#hunterbidenslaptop' hashtag that was picked up in the Biden07-06 data frame, the Ukraine conspiracy theory has taken a similar decrease in popularity. As this particular conspiracy was likely

Results from the Trump-Harris data frame are rather random, but that can be explained by the fact that it is far too small to draw any reasonable conclusions, coming in at only 555 tweets.

In fact, all data frames besides Trump, Biden, and Trump-Biden were rather thin, all coming in below 15,000 entries. So I decided for the 2022 data to focus on the Trump, Biden, and Trump-Biden data frames, and only on the QAnon conspiracy within those data frames.

## 5.1.2 Finding New Hashtags

So I set out to see if there was a set of new hashtags that the QAnon conspiracy community was using on Twitter, using the same methodology I did when first generated my lists of conspiracy theory related hashtags.

I attempted to locate a list of any current popular QAnon hashtags but was unable to find anything, so decided to work from my pre-existing list. While not ideal, some of the old hashtags would still be seeing some use and would hopefully provide links to new hashtags.

However, running the old list of hashtags through the Hashtag Extractor proved far less useful than expected.

```
Tags from Initial List: [('#qanon', 199), ('#pizzagate', 72), ('#maga', 54), ('#trumpwon', 23), ('#lgb', 21), ('#letsgobrando
n', 21), ('#greatawakening', 21), ('#gqp', 15), ('#wwg1wga', 11), ('#qanons', 11), ('#linwood', 5), ('#thegreatawakening', 2),
('#thestorm', 2), ('#redpill78', 1), ('#wearethenewsnow', 1), ('#pedogate', 1)]

Extracted Tags: [('#demvoice1', 362), ('#bluevoices', 362), ('#trublue', 362), ('#trump', 122), ('#trumpcult', 52), ('#illumin
ati', 39), ('#qanoncult', 35), ('#chemtrails', 34), ('#gop', 32), ('#soros', 29), ('#qdupes', 29), ('#birthers', 28), ('#truthe
rs', 28), ('#plandemic', 27), ('#biden', 23), ('#factsmatter', 21), ('#antivaxxers', 21), ('#newworldorder', 19), ('#sethrich',
19), ('#benghazi', 19), ('#obamagate', 19), ('#newsupdate', 18), ('#donaldtrump', 18), ('#sharpiegate', 16), ('#smartnews', 1
5), ('#truth', 15), ('#chavez', 15), ('#hunterbiden', 14), ('#democrats', 14), ('#conspiracytheory', 14), ('#lizardpeople', 1
4), ('#gopcoverup', 13), ('#freedom', 13), ('#conspiracytheories', 13), ('#uranium1', 13), ('#thebiglie', 12), ('#4chan', 12),
('#nwo', 12), ('#covid', 12), ('#january6thhearings', 12), ('#cult45', 12), ('#mindcontrol', 12), ('#jfk', 12), ('#45s', 12),
('#flatearth', 11), ('#wakeup', 11), ('#truthseeker', 11), ('#votebluetosavedemocracy', 11), ('#antivaxxer', 11), ('#vaccinechi
ps', 11), ('#goodbyegop', 10), ('#roevwade', 10), ('#scotus', 10), ('#usa', 10), ('#trumpcultists', 10), ('#fascism', 9), ('#bi
glie', 9), ('#epstein', 9), ('#ghislainemaxwellclientlist', 9), ('#wherearethechildren', 9)]
```

**Figure 14. Initial List of Hashtags (top) Extract Hashtags from 2022 Data (bottom)**

There is a very small number of hashtags from the existing conspiracy hashtag list that have appeared across the 2022 data. Comparing this to results from the 2020 data and accounting for the difference in size, this constitutes a %%% decrease in QAnon hashtag use.

What's more, a large number of these extracted hashtags have no relation to QAnon. The tags '#demvoice1', '#bluevoices', and '#trueblue', are from one tweet that received several retweets that merely mentioned QAnon. Then a great deal of what may be expected to be QAnon or at least conspiracy related such as '#truthers', '#sharpiegate, '#soros', '#deepstate', '#lizardpeople', is again from a

single retweet calling trump cultists deluded, followed by a mocking string of 18 ridiculous conspiracy tags.

Some of the remaining tags such as '#gqp' (a portmanteau of GOP and Q, referring to Republican politicians supporting QAnon conspiracy theories) and '#trumpwon' imply QAnon membership. I also opened up to more generic hashtags that may not have QAnon implications, such as '#maga','#trumpwon','#letsgobrandon', and '#lgb'.

Despite adding these to the list of hashtags and running it through the Hashtag Extractor again, it returned no results that could be linked to QAnon communities.

It is quite likely that there are bubbles of QAnon communities that were either hiding in the data or that were simply not picked up in the data gathering. However, possibly a better explanation is that QAnon conversation has simply left Twitter in favour of other social media platforms such as Parler and Truth Social, which are more accepting of those kinds of beliefs.

## 6.1 Conclusions

The goals of this project as set were rather broad, and while there is much refinement that could be done and further research to be completed, I believe they were all met.

Firstly, through the use of various python libraries such as os and gzip I was able to handle the data within the huge misinformation callout corpus, and organise it into many reasonably sized data frames with pandas. From these data frames, I could easily select the appropriate data frame/s for the problem at hand with the powerful tools afforded me by pandas.

Within the data itself, there were several things found that distinguished the presence and prevalence of conspiracy theories and specific sentiments between the various figures and pairings. Perhaps the largest, and the most unsurprising, was that of the Trump data frame. It had the highest percentage of Negative sentiment, as well as the highest percentage of conspiracy related tweets. This carried over to any data frame which was another figure who was paired with Trump, as his inclusion in a pairing caused both negative sentiment and conspiracy prevalence to rise in comparison to the original figure's solo data frame.

I became acquainted with the VADER sentiment analysis tool and found it to be a very useful tool. Part of what makes it so useful is not just the gold-standard lexicon that the creators have laboriously created, but also the ease with which somebody could make additions to the lexicon to make it fit their purpose.

## 6.2 Future Work

For future work, there is a wide breadth of directions I could envision taking this project just within the limits of the approaches taken over the course of this project.

It would be interesting to expand the scope of the data explored. An earlier starting point would cover the Democratic Primaries and could benefit from looking at other prominent figures, the controversial figure of Bernie Sanders would undoubtedly reveal some things of interest.

A later endpoint could stretch up to the election, which could be very interesting to see the discourse grow as Trump's defeat loomed and Biden was finally victorious. Qanon activity would probably grow, and while the dataset used in this project did not have much to go on regarding accusations of election fraud, that kind of discourse was very common leading up to election day and the immediate aftermath.

An even later date, which could be the most interesting of all would have been to look at the January 6[th] Insurrection, as this was the point where followers of the qanon conspiracy theory truly came into the limelight. The immediate aftermath would also be interesting to see, as there was a major culling of qanon Twitter

accounts after the fact (BBC News 2021), so it could be researched how effective that was in curtailing qanon discourse on the platform.

Another time period to be explored would be the 2016 election, and see how the prevalence of conspiracy theories between the two differed or were similar. The pizzagate conspiracy that existed during the 2016 election was a precursor to qanon after all.

Other social media platforms could be interesting to explore. While places like Reddit and 4chan (or other imageboards) have proven to be havens for conspiracy theorists, platforms such as Parler and Truth Social would be very interesting to look at.  These are platforms that the alt-right have adopted after their exile from Twitter, with the latter being founded by Trump himself.

In terms of technical improvements that could be made that I was unable to pursue in the scope of this project, further exploration of sentiment analysis would be interesting. One avenue to pursue would be adding to the VADER lexicon to improve the analysis of conspiracy theories, or potentially exploring other forms of sentiment analysis.

Another interesting area that I would have liked to try, but did not have the proper time to cover in the appropriate depth, is that of machine learning. It can be applied to NLP, with research done showing ML models being very effective at identifying conspiracy theories (Marcinello et al 2021).


## 6.3 Reflection

The process of developing this project has been challenging, made none the easier that I chose to engage in an area of study with which I had very little experience previously.

Throughout the course of this project, I have learned a great deal having greatly developed my Python skills as well as coming to grips with the exciting field of Natural Language Processing and dipping my toes into Machine Learning.

I've also become well acquainted with the pandas library, and am confident in creating and manipulating data frames. I also had to learn to manage and explore huge quantities of data.

However, I feel the organisation or structure of the development of my project could have benefitted from more forethought and planning. While the topic of the 2020 US election was of interest to me, I felt rather unsure of the direction to take the overall project.

I think if I had established a more concrete objective earlier in the research process the overall result of my project could have greatly benefitted. I felt drawn in many directions with my research, and as such do not feel I took a deep enough study of

any particular area. Still, failures in the organisation and planning a project taught me important lessons in self-management in projects of this nature for the future.

**References**

1. BBC News. 2022. *Twitter suspends 70,000 accounts linked to QAnon*. [online] Available at: https://www.bbc.co.uk/news/technology-55638558 [Accessed 22 September 2022].

2. Bird, S., Klein, E. and Loper, E., 2009. *Natural language processing with Python*. Beijing: O'Reilly. Liddy, E.D. 2001. Natural Language Processing. In Encyclopedia of Library and Information Science, 2nd Ed. NY. Marcel Decker, Inc.

3. Chowdhary, K.R. (2020). Natural Language Processing. In: Fundamentals of Artificial Intelligence. Springer, New Delhi. https://doi.org/10.1007/978-81-322-3972-7_19

4. Dilley, L., Welna, W. & Foster, F. (2021/2022). QAnon Propaganda on Twitter as Information Warfare: Influencers, Networks, and Narratives. Accepted Sept. 16, 2021 at Frontiers in Communication, 6:707595, doi: 10.3389/fcomm.2021.707595

5. Facebook. 2022. *2022 Q2 Report*. [online] Available at: https://s21.q4cdn.com/399680738/files/doc_financials/2022/q2/Meta-06.30.2022-Exhibit-99.1-Final.pdf [Accessed 20 September 2022].

6. Giachanou, A. and Crestani, F., 2016. Like It or Not. *ACM Computing Surveys*, 49(2), pp.1-41.

7. Godfrey et al. (2014). A Case Study in Text Mining: Interpreting Twitter Data From World Cup Tweets.

8. Hunter, J., 2007. Matplotlib: A 2D Graphics Environment. *Computing in Science & Engineering*, 9(3), pp.90-95.

9. Hutto, C., n.d. *GitHub VADER Sentiment Analysis.* [online] GitHub. Available at: https://github.com/cjhutto/vaderSentiment [Accessed 22 September 2022].

10. Hutto, C. and Gilbert, E. (2014) VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text. *Proceedings of the International AAAI Conference on Web and Social Media*, 8(1), pp. 216-225. Available at: https://ojs.aaai.org/index.php/ICWSM/article/view/14550 [Accessed: 22 September 2022]

11. Institute for Strategic Dialogue. 2020. *The Genesis of a Conspiracy Theory.* Available at: https://www.isdglobal.org/wp-content/uploads/2020/07/The-Genesis-of-a-Conspiracy-Theory.pdf

12. Khan, G., Swar, B. and Lee, S., 2014. Social Media Risks and Benefits. *Social Science Computer Review*, 32(5), pp.606-627.

13. Mahase, E., 2020. Hydroxychloroquine for covid-19: the end of the line?. *BMJ*, p.m2378.

14. Marcellino, William, Todd C. Helmus, Joshua Kerrigan, Hilary Reininger, Rouslan I. Karimov, and Rebecca Ann Lawrenc. 2020, *Detecting Conspiracy Theories on Social Media: Improving Machine Learning to Detect and Understand*. Online Conspiracy Theories, Santa Monica, Calif.: RAND Corporation, RR-A676-1, 2021. [Accessed 22, 2022] Available at: https://www.rand.org/pubs/research_reports/RRA676-1.html

15. Marwick, A. and Lewis, R. 2017. *Media Manipulation and Disinformation Online.* Data & Society, pp35-36.

16. Mirani, T.B., & Sasi, S. (2016). *Sentiment Analysis of ISIS Related Tweets Using Absolute Location.* 2016 International Conference on Computational Science and Computational Intelligence (CSCI), 1140-1145.

17. NPR.org. 2022. *An attempted attack on an FBI office raises concerns about violent far-right rhetoric*. [online] Available at: https://www.npr.org/2022/08/12/1117275044/an-attempted-attack-on-an-fbi-office-raises-concerns-about-violent-far-right-rhe [Accessed 22 September 2022].

18. Ordun, Catherine & Purushotham, Sanjay & Raff, Edward. (2020). Exploratory Analysis of Covid-19 Tweets using Topic Modeling, UMAP, and DiGraphs. Available at: https://arxiv.org/pdf/2005.03082.pdf Accessed 22 September 2022]

19. The pandas development team, 2020. *pandas-dev/pandas: Pandas*.

20. Preece, A., Spasic, I., Evans, K., Rogers, D., Webberley, W., Roberts, C. and Innes, M., 2018. Sentinel: A Codesigned Platform for Semantic Enrichment of Social Media Streams. *IEEE Transactions on Computational Social Systems*, 5(1), pp.118-131.

21. Roesslein, J., 2020. Tweepy: Twitter for Python! URL: *https://github.com/tweepy/tweepy.*

22. Shao, C., Ciampaglia, G., Varol, O., Yang, K., Flammini, A. and Menczer, F., 2018. The spread of low-credibility content by social bots. *Nature Communications*, 9(1).

23. Stecula, D. and Pickup, M., 2021. Social Media, Cognitive Reflection, and Conspiracy Beliefs. *Frontiers in Political Science*, 3.

24. Twitter, 2022. *Twitter 2022 Q2 Report*. [online] S22.q4cdn.com. Available at: https://s22.q4cdn.com/826641620/files/doc_financials/2022/q2/Final_Q2'22_Earnings_Release.pdf [Accessed 22 September 2022].

25. Wardle, C. and Derakhshan, H. 2017. *Information disorder: Toward an interdisciplinary framework for research and policymaking.* Pp. 21-22, 49-56. Available at: http://tverezo.info/wp-content/uploads/2017/11/PREMS-162317-GBR-2018-Report-desinformation-A4-BAT.pdf [Accessed 22 September 2022].

26. Yaqub, U., Chun, S., Atluri, V. and Vaidya, J., 2017. Analysis of political discourse on twitter in the context of the 2016 US presidential elections. *Government Information Quarterly*, 34(4), pp.613-626.