s s	Let's see an example of the top posts for both users: Shittymorph For those wondering Banksy was either in the auction house himself or had a helper present that likely shredded the artwork with a remote control. News reports are saying that "a man dressed in black sportic sunglasses and a hat was seen scuffling with security guards near the entrance to Sotheby's shortly after the incident." It wasn't immediately clear if the man had been detained. The auction house believes the shredding mechanism was activated remotely and powered by batteries that had been built into the frame back in nineteen ninety eight when the undertaker threw mankind off hell in a cell, and plumments feet through an announcer's table.
T	CuyWithRealFacts This phenomenon happens when a desert floods, and the water flows over particularly dry, deep sand. The water will seep into the ground so fast that it parts the sand and causes a clear channel down to an underground lake. These subterranean desert lakes sit so deep below the surface that we've never clearly seen them. Drilling to them is impossible in the sand, so they largely remain a mystery. The closest we've gotten is by sand the sand, so they largely remain a mystery.
te T re h	tethered cameras into the whirlpools as they form for some grainy footage before the water channel closes. The subterranean lakes seem to stretch for miles in every direction, but the cameras are always removed from their tether and carried off by unseen creatures soon after they arrive. The cameras have been recovered months later - wrapped in presents under Christmas trees from Santa Claus or in Easter Baskets - leading to the conclusion that all of the magical creatures on Earth vacation at those antipos://www.reddit.com/r/gifs/comments/ayjmgv/terrifying_whirlpools_in_the_las_vegas_desert/ei1gp83/ The bold shows the distinct ending of the post for each user. Some have speculated if both accounts are operated by the same person.
	2) Aims and objectives will be exploring the follow in this project: • Analyze user habits of both accounts • Measure and compare both accounts' operating habits • Compare the frequently visited subreddits
- *;	 Analyze comment content Which words do the accounts prefer? Analyze the comments in other quantitative measurements Look for a relationship or common element between the data Analyze unique writing methods
for E	Data gathering method Datasets for both accounts were acquired through the use of PRAW: The Python Reddit API Wrapper. This API allowed me to gather the posts from both accounts and also other related data and export it to format for further processing. Since shittymorph only had a post count of ~370, I decided to limit the dataset for GuyWithRealFacts to 400 top entries as to create a more equal representation. Data from the specific comments from the account was organized into rows and columns with headers that correspond to the comment and its attributes. This way of organizing data will make it easier for further analysis when comparing attributes.
E C	Choice of user account Both accounts are very similar in that the user posts some comment that leads the reader to believe the expertise of the speaker. Later on, a text based rickroll is introduced for fun. Data suitability The appropriateness and usability of this method of data gathering is rated very highly. Accuracy is high as the API grabs the actual data of the comments and its attributes in real time. Usability is high as data is organized into rows and columns. The PRAW API also allowed me to write to CSV and specify the layout of the data. This is very important to have a streamlined data processing parts.
E	The reason for selecting the two accounts were because both are novelty accounts, and both operate in the same domain. That is, they both create comments for the sole purpose of luring the reader in with elevated storytelling skills and then hit the reader in the end with some signature move. Shittymorph was first selected as I have personally been involved in reading of the account's posts. A suitable match wo ocated via Google search and that led me to wonder if both accounts are operated by the same person. Ethics of gathered data Data on the comment and its attributes were gathered through the PRAW API, and as far as I can tell, all the content solely belongs to the user and reddit as per the TOS[2]. Content submitted to the website automatically available for reddit to use in their operations.
	import praw import csv import os reddit = praw.Reddit(client_id="redacted",
	client_secret="redacted", user_agent="my user agent",) Create CSV file and write the data from PRAW API into it # Open a CSV file for write access (overwrites existing file) with open('filename1a.csv', 'w', newline='', encoding='utf-8-sig') as csv_file:
	<pre># Get a CSV writer object. The 'excel' dialect works well for most things. csv_writer = csv.writer(csv_file, dialect='excel') # Write CSV headers csv_writer.writerow(['id', 'post-created-at', 'text', 'title', 'comment-subreddit', 'comment-subreddit-subscriber', 'comment-submission-subreddit', 'upvotes', # Iterate through prettyoaktree's comments. # If you want to limit the number of comments returned, specify the number in the limit argument. Can change limit to 370, or more for comment in reddit.redditor('shittymorph').comments.top(limit=1): # Write both the comment's text and the parent post's title to the csv file</pre>
:	<pre>csv_writer.writerow([comment.id, comment.created_utc, comment.body, comment.submission.title, comment.subreddit, comment.subreddit.subscribers, comment.su # All purpose import numpy as np import pandas as pd import string # Visualisation import matplotlib.pyplot as plt import matplotlib.pyplot as plt</pre>
	<pre>import seaborn as sns import matplotlib from matplotlib import rcParams # Regular expression library import re # Wordclouds from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator from PIL import Image</pre>
	<pre># Time manipulation from datetime import datetime import pytz from pytz import timezone # nltk import nltk from nltk.tokenize import word_tokenize from nltk.corpus import stopwords from nltk.probability import FreqDist</pre>
:	<pre>from nltk.tokenize import RegexpTokenizer from nltk.stem import WordNetLemmatizer # load and read the dataset df_sm = pd.read_csv('sm-top-400-sample.csv', encoding='utf-8-sig') df_gwf = pd.read_csv('gwf-top-400-sample.csv', encoding='utf-8-sig') # convert unix time to datetime object df_sm['post-created-at'] = pd.to_datetime(df_sm['post-created-at'], unit='s') df_sm['post-created-at'] = pd.to_datetime(df_sm['post-created-at'], unit='s')</pre>
:	df_gwf['post-created-at'] = pd.to_datetime(df_gwf['post-created-at'], unit='s') id
ı	elf81 j 2016-06-28 od:30:17 undertaker mankind thin today, in nineteen ninety eig videos 26493919 videos 45424 3373 0.78 1.530159e+09 /r/videos/comments/8ug87l/twenty_years_age that the data from PRAW API allowed us to extract related information from every commemnt from both users. Dataset columns: Id: Unique ID assigned to every comment
	 post-created-at: Comment creation time and date in GMT text: The body of the comment, as Markdown. title: The title of the submission. comment-subreddit: Name of subreddit the comment belongs to comment-subreddit-subscriber: Number of subscribers to subreddit comment-submission-subreddit: Name of subreddit the comment belongs to upvotes: Number of total upvotes for comment numb-comments: Total number of comments under original submission
F	 upvote-ratio: Ratio of total upvotes for submission submission-created-at: Submission creation time permalink: Link to comment For the first comparison, let's do a simple comparison on the creation times of the comment. We can further understand the daily schedule and user through the posting habits. For abbreviation purposes in the writeup, we will refer to shittymorth as ("SM"), and GuyWithRealFacts as ("GWF").
:	<pre># timeshift from UTC time to Pacific time TIMESHIFT = -7 # count posts for each hour for sm and gwf posts_by_hour_sm = df_sm['post-created-at'].dt.hour.value_counts().sort_index() posts_by_hour_gwf = df_gwf['post-created-at'].dt.hour.value_counts().sort_index() # plot posts by sm ax = sns.barplot(x=(posts_by_hour_sm.index + 24 + TIMESHIFT) % 24,</pre>
	<pre>plt.xlabel('hour') plt.ylabel('count of posts') plt.show() #plot the posts by GWF ax2 = sns.barplot(x=(posts_by_hour_gwf.index + 24 + TIMESHIFT) % 24,</pre>
	plt . show() shittymorph posts by hour (Pacific GMT -7) 25 - 20 - 20 - 20 - 20 - 20 - 20 - 20 -
	GuyWithRealFacts posts by hour (Pacific GMT -7)
	The state of the s
	Analysis of above data: Some time zones were trialed, but the (Pacific GMT -7) time seems to fit the schedule of both users.
ir T	We can see that SM prefers to post during the evening, with the peak coinciding with after dinner time. The trend appears to be less likely to post during the morning and noon, and more likely with chances ncreasing throughout the day all the way until around 21:00, where it starts to decrease. The data from GWF shows a different story. It seems GWF prefers to post during 2 distinct peaks, with the first peak in the morning, and the second in the evening around dinner time. There seems to be a disposting chance during the 11:00 to 16:00. What we can infer from this is that GWF seems to have a daily schedule that he follows and perhaps this is a sign that GWF has a stable day job or is attending schedule. 1. Not many posts overnight 2. Not many posts during midday after lunch 3. Peak posting is before and after work time.
- L	Let's plot the number of posts both accounts make per day. from statistics import mean day = df_sm.groupby(by=df_sm['post-created-at'].dt.date)['text'].count() print("Shittymorph mean value of posts per day:", mean(day.values))
	<pre>ax = sns.barplot(x=day.index, y=day.values) ax.set_title('Shittymorph number of posts per day:') # rotate labels ax.set_xticklabels(ax.get_xticklabels(), rotation=30) nbins = 8 # show only nbins labels 1 = ax.xaxis.set_major_locator(matplotlib.ticker.MaxNLocator(nbins=nbins))</pre> Shittymorph mean value of posts per day: 1
	Shittymorph number of posts per day: 8-7-6-5-4-3-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1
	2- 1- 2016-07-27 2017-02-29 2017-07-3018-02-14 2019-02-20 2012-08-05 post-created-at
	<pre>## plot gwf day = df_gwf.groupby(by=df_gwf['post-created-at'].dt.date)['text'].count() print("GuyWithRealFacts mean value of posts per day:",mean(day.values)) ax = sns.barplot(x=day.index, y=day.values) ax.set_title('GuyWithRealFacts number of posts per day:') # rotate labels ax.set_xticklabels(ax.get_xticklabels(), rotation=30) nbins = 8 # show only nbins labels</pre>
	1 = ax.xaxis.set_major_locator(matplotlib.ticker.MaxNLocator(nbins=nbins)) GuyWithRealFacts mean value of posts per day: GuyWithRealFacts number of posts per day: 10
	2 2 2 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
E T	Analysis of above data: Both accounts show that usually only one comment is made per day. 2017 shows some interesting activity for both accounts as both accounts deviated from the normal practice and posted multiple times per There is no overlap in the peak for both as GWF's peak is in July 2017, and SM already settled into its 1 post per day habit. Let's now plot the preferred subreddits of both users.
: F	#plot for shittymorph print("shittymorph upvotes") sns.catplot(x="comment-subreddit", y="upvotes", data=df_sm, aspect=18/9) plt.xticks(rotation=90) # rotate the x-axis labels plt.show() #plot for gwf print("GuyWithRealFacts upvotes")
	<pre>sns.catplot(x="comment-subreddit", y="upvotes", data=df_gwf, aspect=18/9) plt.xticks(rotation=90) # rotate the x-axis labels plt.show() shittymorph upvotes 40000 -</pre>
	30000 - 20000 - 10000
	o decided interesting announcements announcements announcements applications interesting astruct worldnews worldnews worldnews worldnews aww worldnews aww ondopoles interesting astruct worldnews aww worldnews aww ondopoles westworld funny funny aports westworld funny funny aports worldnews worldnews worldnews worldnews aww on the property of the pr
	14000 - 12000 - 10000 - 10000 -
	aww
F	Analysis of the above data:
T T -	The story is slightly different form GWF. It seems this account yielded the most upvotes from gifs, aww, interesetingasfuck, and pics subreddits. The only similarity is that posts for both users in the gifs subreddit yield a similar spread of upvotes, with SM reaching the higher limit of what is possible. Next, we will plot the number words in each comment made by both users. Perhaps, this can help us visualize any overlap in writing style or preferences between the two users.
: [First we will create a new colum "word_count", then use the str.split str.len function to get the total words in all the comments in our datafame df_sm['word_count'] = df_sm['text'].str.split().str.len() df_gwf['word_count'] = df_gwf['text'].str.split().str.len() # import warnings filter from warnings import simplefilter # ignore all future warnings
	<pre>simplefilter(action='ignore', category=FutureWarning) ax = plt.gca() ax.get_xaxis().get_major_formatter().set_scientific(False) sns.distplot(df_sm["word_count"], label='shittymorph') sns.distplot(df_gwf["word_count"], label='GuyWithRealFacts') plt.legend() plt.show()</pre>
	0.012 -
	Analysis of the above data: We can see that SM prefers short comments, with 2 peaks at around under 50 and at 100 words.
T _ L	GWF's comment show a different story. WE can see that this user prefers longer comments, with the distribution peak of 200. The bulk of the word count of the comments can be interpreted as 100 – 300. This shows very little similarity between both users. Let's dig deeper into the data and create a word cloud for the common words in our dataset. This can help us understand the vocabulary and topics of both users. First, we will take the text in the dataframe and tokenize it and assign it to tokenized_sm and tokenized_gwf. Punctuation will also be removed in this stage.
n V	Because SM has a signature ending has the same copypasta every time: nineteen ninety eight when the undertaker threw mankind off hell in a cell, and plummeted sixteen feet through an announcer's table. We will use the str.replace and give it 2 arguments. First is the beginning of the copypasta which is either 1998 or nineteen, then we remove everything between that word and the last—which is table. Data cleaning and processing
	<pre>tokenize_sm = df_sm['text'] ## remove copypasta from end # remove puncuation tokenize_sm = tokenize_sm.str.replace(r'[^0-9a-zA-Z\s]+', '', regex=True) for x in tokenize_sm: #remove all text between "nineteen" and "table", which is the signature of sm tokenize_sm = tokenize_sm.str.replace(r'(?s)nineteen.*?table', '', regex=True) tokenize_sm = tokenize_sm.str.replace(r'(?s)1998.*?table', '', regex=True)</pre>
:	<pre>#tokenize sm dataset tokenize_sm = tokenize_sm.apply(word_tokenize) stop = stopwords.words('english') # remove stopwords tokenize_sm = tokenize_sm.apply(lambda x: [item for item in x if item not in stop])</pre>
	<pre>tokenize_gwf = df_gwf['text'] ## remove puncuation tokenize_gwf = tokenize_gwf.str.replace(r'[^0-9a-zA-Z\s]+', '', regex=True) #tokenize gwf dataset tokenize_gwf = tokenize_gwf.apply(word_tokenize) stop = stopwords.words('english') # remove stopwords tokenize gwf = tokenize gwf.apply(lambda xy. [item for item in x if item not in stop])</pre>
: [<pre>tokenize_gwf = tokenize_gwf.apply(lambda x: [item for item in x if item not in stop]) def lemmatize_text(text): lemmatizer = WordNetLemmatizer() return [lemmatizer.lemmatize(w) for w in text] tokenize_sm.apply(lemmatize_text) tokenize_gwf.apply(lemmatize_text) 0 [This, phenomenon, happens, desert, flood, wat</pre>
	<pre>[Dr, Mrs, Biden, would, least, light, security [For, long, time, early, day, metal, working, [The, whole, movement, called, indenting, left [The, best, thing, cat, lot, good, thing, pick [Throughout, time, lot, Canadas, wildlife, spo [This, phenomenon, actually, actually, documen [In, theory, hot, sauce, hot, enough, acid, ge [The, dad, reflex, actually, innate, gained, f [Some, dork, made, post, seemed, real, first,</pre>
	<pre>Name: text, Length: 370, dtype: object # plot wordcloud for sm print("shittymorph wordcloud") text = tokenize_sm.values # Generate a word cloud image wordcloud = WordCloud(background_color="white").generate(str(text)) # Display the generated image: plt.figure(figsize=(10,5))</pre>
	<pre>plt.imshow(wordcloud, interpolation='bilinear') plt.axis("off") plt.show() # plot wordcould for gfw print("GuyWithRealFacts wordcloud") text2 = tokenize_gwf.values # Generate a word cloud image for gwf wordcloud = WordCloud(background_color="white").generate(str(text2))</pre>
	# Display the generated image: plt.figure(figsize=(10,5)) plt.imshow(wordcloud, interpolation='bilinear') plt.axis("off") plt.show() shittymorph wordcloud Informatic tida train
	could house time back if the work was a state of the work was take was a state of the work was a state
	Idistract fact that around work inever thats in the state of the state
	known' due instead
	The state of the s
J III	Analysis of the above data: It seems like the vocabulary and word choice is quite different between the two users. SM prefers to use "I" quite a bit and that suggests that the user writes in the first-person tone.
H III S	Analysis of the above data: It seems like the vocabulary and word choice is quite different between the two users. SM prefers to use "I" quite a bit and that suggests that the user writes in the first-person tone. When looking at GWF's wordclould, it becomes more interesting. GFW seems to like to use the word cat, cats, and dogs quite a bit. This may suggest that GFW is a cat owner or perhaps a cat lover. When referencing against the common subreddits GFW visits, it makes more sense as "aww" is at number 2. Visitors to that subreddit can mostly see cute animal media. Summary Conclusions Through this data analysis on the posts made by both accounts, it is noted that many Metris and results of the analysis does not lead to the same person operating SM and GWF.
# HH S V rr - S C C T C C C C C C C C C C C C C C C C	Analysis of the above data: It seems like the vocabulary and word choice is quite different between the two users. SM prefers to use "I" quite a bit and that suggests that the user writes in the first-person tone. When looking at GWF's wordclould, it becomes more interesting. GFW seems to like to use the word cat, cats, and dogs quite a bit. This may suggest that GFW is a cat owner or perhaps a cat lover. When eferencing against the common subreddits GFW visits, it makes more sense as "aww" is at number 2. Visitors to that subreddit can mostly see cute animal media. Summary Conclusions Through this data analysis on the posts made by both accounts, it is noted that many Metris and results of the analysis does not lead to the same person operating SM and GWF. would say the most telling sign is the post creation time as both shows that the accounts owned by people living different lifestyles. The wordclould also shows that the vocabulary choice of both accounts and assimilar. References: 1) https://reddit.zendesk.com/hc/en-us/articles/205313845-What-do-these-expressions-mean-22 https://www.redditirc.com/policies/user-agreement
HIII S V rr T I I d F F N F H	Analysis of the above data: It seems like the vocabulary and word choice is quite different between the two users. SM prefers to use "I" quite a bit and that suggests that the user writes in the first-person tone. When looking at GWF's wordclould, it becomes more interesting. GFW seems to like to use the word cat, cats, and dogs quite a bit. This may suggest that GFW is a cat owner or perhaps a cat lover. When referencing against the common subreddits GFW visits, it makes more sense as "aww" is at number 2. Visitors to that subreddit can mostly see cute animal media. Summary Conclusions Through this data analysis on the posts made by both accounts, it is noted that many Metris and results of the analysis does not lead to the same person operating SM and GWF. would say the most telling sign is the post creation time as both shows that the accounts owned by people living different lifestyles. The wordclould also shows that the vocabulary choice of both accounts are dissimilar. References: 1) https://reddit.zendesk.com/hc/en-us/articles/205313845-What-do-these-expressions-mean-