

YouNICON: YouTube’s CommuNity of Conspiracy Videos

Liaw Shao Yi¹, Fan Huang², Fabricio Benevenuto³, Haewoon Kwak², Jisun An²

¹School of Computing and Information Systems, Singapore Management University, Singapore

²Luddy School of Informatics, Computing, and Engineering, Indiana University Bloomington, USA

³Federal University of Minas Gerais Belo Horizonte, Brazil

shaoyi.liaw.2022@phdcs.smu.edu.sg, huangfan@acm.org, fabricio@dcc.ufmg.br, haewoon@acm.org, jisunan@acm.org

Abstract

Conspiracy theories are widely propagated on social media. Among various social media services, YouTube is one of the most influential sources of news and entertainment. This paper seeks to develop a dataset, YOUNICON, to enable researchers to perform conspiracy theory detection as well as classification of videos with conspiracy theories into different topics. YOUNICON is a dataset with a large collection of videos from suspicious channels that were identified to contain conspiracy theories in a previous study (Ledwich and Zaitsev 2020). Overall, YOUNICON will enable researchers to study trends in conspiracy theories and understand how individuals can interact with the conspiracy theory producing community or channel. Our data is available at: <https://doi.org/10.5281/zenodo.7466262>.

Introduction

Conspiracy theories are nothing new in human history. Scholarly research on conspiracy theories began in the 1930s (Butter and Knight 2018), and it has been a field that is highly multidisciplinary and diverse (Mahl, Schäfer, and Zeng 2022). Various researchers have proposed definitions for conspiracy theories. Keeley (1999) defines conspiracy theory as “a proposed explanation of some historical event (or events) in terms of the significant causal agency of a relatively small group of persons—the conspirators—acting in secret.” A more general definition of conspiracy theory is provided by Wood, Douglas, and Sutton (2012) as “a proposed plot by powerful people or organizations working together in secret to accomplish some (usually sinister) goal.”

Conspiracy is sometimes considered a form of misinformation. Misinformation is commonly defined as “false or inaccurate information that ... spread regardless of an intention to deceive. (Tomlein et al. 2021)” This suggests that any malicious intent of the content creator is not a necessary condition of misinformation, but incorrectness of information is. Thus, there is a stark difference between conspiracy and misinformation; the intent of powerful people (or organizations) is crucial for the definition of conspiracy. This difference proves the need for in-depth research on conspiracies that should be differentiated from those on misinformation.

A belief in conspiracy often correlates with anomia, lack of interpersonal trust, and having political beliefs at extreme ends of the political spectrum (especially on the right-hand extreme) (Goertzel 1994; Sutton and Douglas 2020). Conspiracy theories, in contrast to non-conspiracy views, tend to be more attractive as they satisfy one’s epistemic (e.g., the desire for understanding, accuracy, and subjective certainty), existential (e.g., the desire for control and security), and social desires (e.g., the desire to maintain a positive image of the self or group) (Douglas, Sutton, and Cichočka 2017). This results in undesirable outcomes like decreased institutional trust and social engagement, political disengagement, prejudice, environmental inaction, and an increased tendency towards everyday crime (Pummerer et al. 2022; Douglas, Sutton, and Cichočka 2017; Jolley et al. 2019). Additionally, conspiracy theories can form a worldview in which believers of a type of conspiracy tend to approve of other conspiracies as well (Wood, Douglas, and Sutton 2012; Dagnall et al. 2015). Polls have also shown that “everyone believes in at least one or a few conspiracy theories (Uscinski 2020)”. Hence, a holistic understanding of conspiracy theories cannot be achieved in isolation from a specific type of conspiracy.

In today’s context, conspiracy theories are widely propagated on social media (Enders et al. 2021; Mahl, Zeng, and Schäfer 2021). Conspiracy narratives are nourished by information cascades on social media and reach a larger audience (Monaci 2021). Consequently, these false narratives tend to outperform real news in terms of popularity and audience engagement within online environments (Coninck et al. 2021; Vosoughi, Roy, and Aral 2018). Enders et al. (2021) show that usage of 4chan/8kun has the highest correlation with the number of conspiracy beliefs, followed by Reddit, Twitter, and YouTube.

Among the social media services, YouTube is one of the most influential sources of news and entertainment (Center 2012). It has 2,562 million monthly active users, and it is the second most popular social network worldwide as of January 2022 (Statista 2022), contributing to a billion hours of video viewed daily (Goodrow 2017). Audit studies show that video recommendations on YouTube can lead to the formation of filter bubbles on misinformation topics (Hussein, Juneja, and Mitra 2020). Similarly, exposure to conspiracy videos might result in undesirable outcomes. For example,

the belief that the 5G cellular network caused COVID-19 has resulted in more than 200 reports of attacks against telecom workers in the United Kingdom (Vincent 2020). The belief in white genocide conspiracies resulted in the death of 51 individuals in New Zealand (Commission et al. 2020). The belief in conspiracy theories is no doubt an issue of concern. However, most existing research focuses only on specific types of conspiracy theories, and not all datasets are available for research communities.

In this work, we build YOUNICON, a curated dataset of YouTube videos from channels identified as producing conspiracy content by Recfluence (Ledwich and Zaitsev 2020). We aim to help researchers to study the patterns of production and consumption of conspiracy videos, such as how individuals interact with those videos from an aggregated (video) or individual level (comments)¹.

YOUNICON comprises the following information:

- Metadata of all 596,967 videos from 1,912 channels that produced conspiracy identified by Recfluence (Ledwich and Zaitsev 2020)
- A list of 3,161 videos manually labeled as being about conspiracy or not
- 37,199,252 comment IDs of comments in all videos with basic metadata and scores from the Perspective API¹
- 100 videos manually labeled for the type of conspiracy.

YOUNICON will be a valuable resource for studying YouTube as a medium of conspiracy theory production and consumption. The contributions of this paper are as follows:

- Curate a large-scale dataset of videos with conspiracy content (<https://doi.org/10.5281/zenodo.7466262>)
- Perform exploratory analyses on the dataset to understand its key properties
- Discuss potential uses for the dataset

Related Work and Datasets

Conspiracy Detection

Table 1 highlights several existing datasets that have been used for conspiracy theory detection research. Existing literature often focuses on misinformation (Lin et al. 2019; Kumar et al. 2020) or specific conspiracy theories related to COVID-19, alien visitation, anti-vaccination, white genocide, climate change, or Jeffery Epstein (Moffitt, King, and Carley 2021; Marcellino et al. 2021; Phillips, Ng, and Carley 2022). Most works focus mainly on Tweets as the unit of the study (Moffitt, King, and Carley 2021; Galende et al. 2022; Phillips, Ng, and Carley 2022; Mahl, Zeng, and Schäfer 2021). For example, Galende et al. (2022) study Tweets explicitly containing the word “conspiracy.” Phillips, Ng, and Carley (2022) have compiled a dataset consisting of four types of conspiracy, namely climate change, COVID-19 origins, COVID-19 vaccine, Epstein Maxwell. Moffitt, King, and Carley (2021) study COVID-19-related conspiracies by training a BERT-based classifier to distinguish conspiracy Tweets.

¹The comment text and real author names are not shared to protect the identity of the commenter. The actual comment text can be rehydrated using YouTube Data API.

Conspiracy Taxonomy

Mahl, Zeng, and Schäfer (2021) used network analysis of co-occurring hashtags in Tweets to assign hashtags into topic groups qualitatively based on their thematic relationship. This resulted in the 10 most visible conspiracies, which include Agenda 21, Anti-Vaccination, Chemtrails, Climate Change Denial, Directed Energy Weapons, Flat Earth, Illuminati, Pizzagate, Reptilians, and 9/11 Conspiracies. While co-occurring patterns of hashtags reveal a partial taxonomy of conspiracy, a more comprehensive one is found on Wikipedia.

On Wikipedia, a list of conspiracy theories is constantly being updated (Wikipedia contributors 2022). Upon closer inspection of the list of conspiracy topics from Wikipedia, we found that it covers well the conspiracies in Mahl, Zeng, and Schäfer (2021) (see Table 3 for details). Hence, we will use the taxonomy of Wikipedia for YOUNICON.

Data Collection

On YouTube, interactions between content creators and consumers occur as follows: a content creator posts a video with a title, description, and tags. A content consumer views, likes, or comments on a video. A “view” represents a playback of a video, a “like” is positive feedback to the video by users, and a “comment” is the way in which online collective debates grow around the video (Bessi et al. 2016; YouTube 2022). A comment can be a reply to a video (a top-level comment) or a reply to other comments.

YouTube Channels about Conspiracy

Ledwich and Zaitsev (2020) curated a list of US-based political channels in the Recfluence project. They classify each channel based on its political leaning, channel type (e.g., mainstream news, AltRight, etc), and topical category (e.g., conspiracy, libertarian, organized religion, LGBT, etc).

We downloaded an entire list of YouTube channels from Recfluence on 25 February 2022 and extracted only channels with the “conspiracy” label. We then used the YouTube Data API² to collect the basic information about these channels. Out of the 2365 channels with the “conspiracy” label, 1912 channels were accessible by the YouTube API. The rest of the channels were deleted from YouTube and thus excluded from the following analysis. While Recfluence provides a quite extensive list of US-based political channels, the resulting list could be improved with more channels. However, all the pipelines used in this work will still be valid.

Video Metadata

In contrast to Recfluence (Ledwich and Zaitsev 2020) that provide channel-level conspiracy information, YOUNICON focuses on video-level conspiracy. For the extracted conspiracy-related channels from Recfluence, we collect the metadata of every video published on those channels. The metadata includes the title, description, tags, number of likes, number of views, duration, and published date. We

²<https://developers.google.com/youtube/v3>

Table 1: Related Datasets for Conspiracy Detection

Dataset	Description	Labels	Annotation Method
Reclufence (Ledwich and Zaitsev 2020)	Political inclination of YouTube channels and tags to characterize channel.	7,085 channels with 2,365 of the channels with the tag conspiracy	Agreement between 3 labelers
Conspiracy theory videos (Faddoul, Chaslot, and Farid 2020)	YouTube Video Dataset of Conspiracy Theory Videos from YouTube’s Watch Next engine	Conspiracy (542) and non-conspiracy (568) in training set	Conspiracy Videos from a Book (Jackson 2017) or Reddit and non-conspiracy videos from random scraping
Conspiracy theory Tweets and videos (Ginossar et al. 2022)	Cross-platform dataset which includes YouTube videos from Tweets related to COVID-19 and vaccines	930,539 Tweets and 1,280 YouTube conspiracy theory videos with transcripts	Not applicable
Twitter conversations conspiracy (Galende et al. 2022)	Twitter conversations dataset of conversations with more than 1000 Tweets that contain the word “conspiracy”.	More than 4,500 conversations	Semi-automatic
Hoaxes and Hidden agendas (Phillips, Ng, and Carley 2022)	Tweets from 4 Topics: climate change, COVID origins, COVID vaccine, and Epstein.	3,100 annotated instances. Conspiracy (2,336) and non-conspiracy (764)	Agreement between 3 labelers.
CMU-MisCOV19 (Memon and Carley 2020)	Tweets related to COVID-19 misinformation with 17 different labels which include topics like Irrelevant, Conspiracy, True Treatment, Politics, Commercial Activity, etc.	4,573 annotated Tweets. Conspiracy (924)	Annotation class determined by 1 single annotator.
Moffitt, King, and Carley (2021)	Extension of the dataset by Memon and Carley (2020) using the same procedure, then collapsed labels into a binary conspiracy classification task.	8,781 labeled Tweets 4,573 Tweets from Memon and Carley (2020) and 4,208 new labeled tweets.	Manually labeled by research assistants.

collect these metadata for 1,049,413 videos in total. To get a better sense of the content that was presented in the videos, we also collect transcripts or subtitles of the videos using a PyPI package, `youtube-transcript-api`³. We only consider those videos with English transcripts, which are 761,565 videos in total.

We further filter out non-English videos by detecting the language of the videos based on their titles, which often summarizes the gist of the video. In particular, we use the Fasttext language identification model (Joulin et al. 2016), which can recognize 176 languages, with a threshold of 0.5 to determine the language with the highest probability for a video. Between the two Fasttext language identification models, we used the larger and more accurate one (i.e., `lid.176.bin`). As a rule of thumb, channels with less than 80 percent of their videos that are in English are excluded from the rest of the analysis.

For all textual metadata, we apply common preprocessing techniques (e.g., remove emojis, URLs, punctuation, and numbers, and convert them to lowercase). Then, we filter out videos that do not have all the metadata. This results in a collection of 596,967 videos with all metadata, which are title,

description, tag, and transcript.

Additionally, we collect top-level comments as a part of the video’s features. We filter out comments’ authors if 1) their comments detected as English are less than 80%, or 2) they leave only one comment. We also eliminate the top-level comments written by the same video creator to focus on the behavior of the viewers. As a result, we obtain 37,199,252 comments. For these comments, we use the Perspective API to perform scoring for toxicity, identity attack, and threat⁴.

Dataset Construction

Figure 1 is a flowchart that summarises the dataset construction proposed in this paper. The following sections will explain the proposed method in detail. Table 2 summarises the variables available in the YOUNICON.

Video Labeling

The procedure of manual annotation is done in accordance with the Institutional Review Board (IRB) guidelines under

³<https://github.com/jdepoix/youtube-transcript-api>

⁴www.perspectiveapi.com

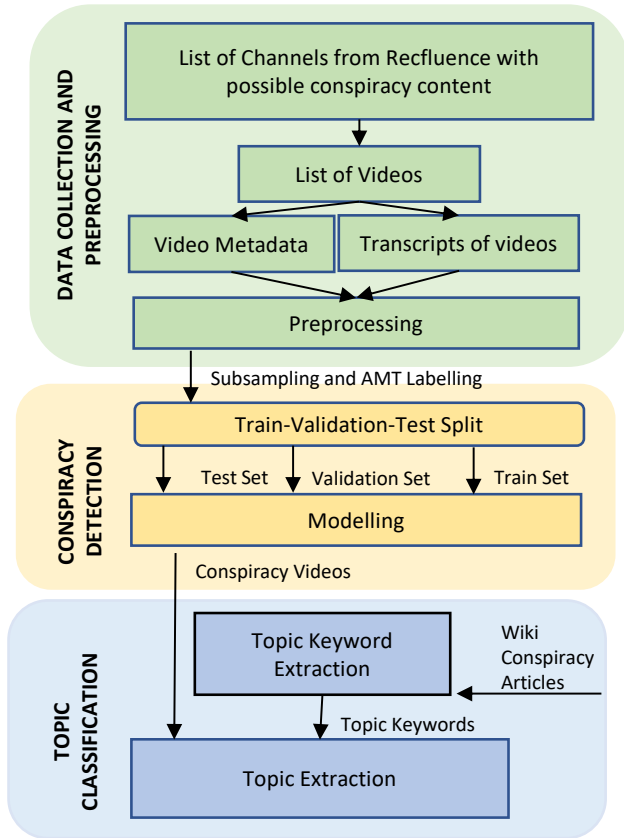


Figure 1: Overview of YOUNICON Construction

the approval number IRB-22-129-A071(922) of Singapore Management University.

We use Amazon Mechanical Turk (AMT) for data labeling. Our labeling task, known as Human Intelligence Task (HIT) in AMT, asks an AMT worker whether a given video contains conspiracy or not. The title, description, tags, and the first 1,000 characters of the transcript of each video are given to AMT workers. We select workers located in the US, with a past HIT approval rate of greater than 98% and 5,000 HITs approved, and compensated them at a rate of 0.05 USD per HIT.

For each video, we recruit three workers and determined a label based on the majority vote. In contrast to misinformation where there is a clear-cut answer, determining whether a video contains conspiracy can be more challenging as an individual’s political or religious belief might affect their decision about conspiracy videos. Thus, we follow a majority voting scheme for each video’s label.

Labeling is conducted in two stages. In the first stage, we sample 2,200 videos. After labeling, we find that this dataset is somewhat imbalanced; only around 20% (436 of the videos out of 2,184) contain conspiracy theories. Although 20% may seem like a relatively large proportion of the videos with conspiracy theories, we note that all these videos are from the channels that are categorized as ‘conspiracy’ in Recfluence (Ledwich and Zaitsev 2020). To make

Table 2: Variables and descriptions in the YOUNICON dataset

Variable	Description
All Videos	
Video_id	YouTube Video ID
Channel_ID	ID of the channel that video is from
Title	Title of the video
Video_Description	Description of Video Provided by content creators
Tags	Tags provided by content creators
Transcript	transcript of the video from the youtube-transcript-api
Published_date	date video is published
Views	number of views
Likes	number of likes
Dislikes	number of dislikes
NumComments	number of comments
Duration	duration
Category	content category provided by YouTube
DurationSec	duration of videos in seconds
Comments	
Comment_Id	ID of comments
VideoId	ID of the video comment is on
Anon_id	anonymised author’s ID
Toxic	Toxic Score from Perspective API
Identity	Identity Attack Score from Perspective API
Threat	Threat Score from Perspective API
LikeCount	number of likes on comment
PublishedAt	publish time of comment
TotalReplyCount	number of replies to comment
Conspiracy Label	
Video_id	YouTube Video ID
Majority_label	1 if the video contains conspiracies 0 otherwise

YOUNICON a better-balanced dataset of conspiracy and non-conspiracy videos, we use Machine Learning models to get pseudo-labels first. We finetune the RoBERTa-large model by using the sampled videos. We split the data into train, validation, and test sets and use the concatenated texts as features for the model. This model attained an accuracy of 0.74, with a positive F1 of 0.5273 and a negative F1 of 0.8207. We used this trained model to assign ‘conspiracy’ or ‘non-conspiracy’ pseudo-labels to all the videos in the full dataset. Then, we sample 1,000 videos with ‘conspiracy’ pseudo-labels and manually labeled them in the same manner. After these 2 rounds of labeling, we obtain a dataset of 3,161 videos (1,144 conspiracy videos (36.2%)). Fleiss’ Kappa, an extension of Cohen’s kappa, is used to measure inter-rater reliability (Fleiss 1971). A score of 0.4111 is calculated, implying that there is a moderate agreement between raters in the dataset (Landis and Koch 1977).

Table 3: Topics of Conspiracy Theories and the corresponding keywords

Topic	Description	Representative Keywords
Aviation	Chemtrails, Air travel and aircraft	chemtrail, black helicopter, airline, aircraft, airlines, remotely, flight, crash
Business and Industry	Deep Water Horizon and New Coca-cola formula	cocacola, deepwater, coke, formula
Deaths and Disappearances	Deaths of prominent leaders and public figures	jfk, dnc, flee, lookalike, assassination
Economics and Society	New World Order, George Soros, Freemasonry	masonic, new world order, george soros, turkey, freemasonry, üst akıl, mastermind, denver airport, economy, freemason, erdoğan, rip
Espionage	Spying with animals or individuals	animal, taliban, spy, wilson, malala, harold, golitsyn
Ethnicity, Race, and Religion	Related to Anti-Religion, Racism, Genicides, and Religious Beliefs	antisemitism, jesus, antichrist, paul, rastafari, catholicism, islamic, bahá'í, bible, apostle, racism, christ, islam, catholic, jihad, bahá'ism, armenianism
Extraterrestrials and UFOs Government, Politics, and Conflict	Alien visitation 9/11. False Flag operations, Political Figures	ufo, anunnaki, extraterrestrial, alien election, congress, trump, marxism, barack, clinton, obama, epstein, ukraine, sandy hook, biden, clintons, national, crisis actor, fema
Medicine	COVID-19 related, Claims that diseases like HIV and Ebola is invented, anti-vaccination, Water fluoridation	vaccination, fluoridation, vaccine, disease, therapy, suppression, virus, pandemic, pharmaceutical, virology
Science and Technology	Global Warming Denial, Flat Earth, Weather Control, Technology Suppression	rfid, weather, earthquake, weaponry, weather control, mkultra, tsunami, haarp, mind control, warming, earth, warm, technology, flat
Outer Space	Staged moon landings by NASA, Nibiru (doomday belief of large planet almost crashing on Earth)	nibiru, outer space, planet, nasa, solar, space

Topic Classification

Going beyond whether a video is about conspiracy or not, we also assign a conspiracy topic to a video based on conspiracy taxonomy compiled on Wikipedia (Wikipedia contributors 2022). In doing so, we first parse the text of the “List of conspiracy theories” page on Wikipedia (Wikipedia contributors 2022). This Wikipedia page contains summaries of the popular conspiracy theories, which include Aviation, Business and Industry, Deaths and Disappearances, Economics and Society, Espionage, Ethnicity, Race and Religion, Extraterrestrials and UFOs, Government, Politics and Conflict, Medicine, Science and Technology, Outer Space and Sports (Table 3). We exclude the category of “Sports” because our dataset, based on Recfluce (Ledwich and Zaitsev 2020), is unlikely to contain Sports-related conspiracies. The topic “Fandom, celebrity relationships, and shipping”, a new topic added on 18 May 2022, which is after our Wikipedia data collection, is also not included in this analysis.

The topic classification consists of two stages: 1) keyword extraction and 2) topic inference. For keyword extraction, we identify representative words of each topic using log-odds ratios with informative Dirichlet priors (Monroe, Colaresi, and Quinn 2008), which is a widely used technique for a large-scale comparative text analysis (An et al.

2021; Kwak, An, and Ahn 2020). It estimates the log-odds ratio of each word between two corpora i and j given the prior frequencies obtained from a background corpus. We rank the words based on their log-odds scores and obtain a list of representative words for each of the conspiracy theories. The background corpus used in this analysis is the “google 1-gram” (Michel et al. 2011), extended with the counts of the vocabulary used in the “list of conspiracy theories” Wikipedia page. For each conspiracy topic, we compare the corpus of one topic against the concatenated corpus of all other topics. For each topic, we use the top ten keywords as the preliminary keywords for topic inference (see Table 3 for the list of words extracted). We also add the subtopic names listed in Wikipedia to the keywords of each topic. We then convert 21 keywords to bigrams or trigrams to be more distinguishable (e.g., the keywords new, world, and order should be considered as a trigram, not three unigrams) and remove 62 keywords related to countries or locations (Malaysia, Wuhan) or those that are generic (January, human).

Having the representative keywords for topics at hand, we infer the topic of the video by simply using a keyword-matching method. We match the keywords in each topic to the video’s features by using spaCy’s PhraseMatcher. We assign a topic by choosing one with the highest frequency of

Table 4: Feature comparisons between conspiracy (C) and non-conspiracy (NC) videos. Avg. and Med. are average and mean values, respectively.

Feature	Avg.(C)	Avg.(NC)	Med.(C)	Med.(NC)
Duration (s)	1,398	1,298	688	638
Likes	736	645	136	72
Comments	171	113	32	14
Views	27,245	20,021	4,317	2,018

keywords in a video’s features.

Exploratory Data Analysis

To provide a brief overview of the dataset, we conduct an exploratory analysis. We first compare the difference in engagement of videos with conspiracy theories and those without conspiracy theories. Table 4 shows the average and median values of various features of conspiracy and non-conspiracy videos. Conspiracy videos have longer lengths and get more likes, comments, and views than non-conspiracy videos. Conspiracy videos have 736 likes, 171 comments, and 27,245 views on average, while non-conspiracy videos have only 645 likes, 113 comments, and 20,021 views. All differences are statistically significant based on the Mann Whitney U test (Mann and Whitney 1947) for unpaired samples.

Results

Conspiracy Detection

We use the annotated data of 3,161 videos to build a classifier that detects whether a video is about conspiracy or not. Since our data is slightly unbalanced (1,144 videos are conspiracy), we perform under-sampling to balance the classes for the training. For testing, we use the holdout test set sampled from the initial (first-round) 2,200 annotated videos.

As a feature, we use all video’s textual meta information, including title, tags, description, and transcript. Since the deep learning models can take 512 tokens at maximum (Liu et al. 2019), we truncate the video description and transcript, using the first 200 tokens. The feature input, called combined, is created by concatenating the first 200 tokens or words for both the video description and transcript, followed by the title and tags.

In order to compare the performance of the models, other than simply accuracy, recall, or precision, we use the F1-score:

$$F_1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

, which is calculated for both the positive and negative classes. To account for class imbalance, F1 weighted, which is the F1 score weighted by the support, is also used.

Table 5 summarizes the prediction results of various models. Dummy Classifier predicts all videos as negative (or non-conspiracy), yielding an accuracy of 0.8, which is the same as the proportion of non-conspiracy videos in the test set. Traditional machine learning models, including Naive

Bayes, Logistics Regression, and Support Vector Machine with Linear Kernel (SVM), are also tested. All three models, Naive Bayes, Logistics, and SVM slightly outperform the Dummy classifier, obtaining a weighted F1 of 0.7141, 0.7930, and 0.7863, respectively.

We also explore pre-trained language models, such as RoBERTa-large. The training set is further split into 80-20 train-validation split for finetuning of the pre-trained model. The learning rate of 1e-5 is used with a batch size of 4 with random seed 13 for finetuning. Our results in Table 5 show that the pre-trained models result in better performance in all metrics but recall, obtaining an accuracy of 0.8575 and weighted F1 of 0.8624.

In Table 5, we also show the prediction results based on individual features. By comparing the weighted F1 of models built based on each feature, we observe that the tags are best among the individual features, followed by video description, titles, and transcript.

Zero Shot and Few Shot Classification

We further conduct experiments to examine if it is possible to detect conspiracy theories via zero and few-shot learning. Zero and few-shot learning are techniques that aim to make predictions for new classes with limited labeled data. We test pre-trained Natural Language Inference (NLI) (Bowman et al. 2015; Williams, Nangia, and Bowman 2018) and Natural Language Generation (NLG) (Lewis et al. 2020; Zhang et al. 2022) models on zero and few-shot settings. We use all the features of videos as the input for those models.

For the NLG models, we test on both auto-regressive generation and sequence to sequence models⁵. However, we find that the generated results of zero-shot and most few-shot models are simply a repeat of the given text, from which we cannot infer the classification labels. The model could generate clear classification indicators (i.e., *yes* or *no* in our setting) only for the few-shot settings with 128 fine-tuning data instances. However, it predicts all inputs as non-conspiracy.

As for the NLI models, we apply the top three most popular fine-tuned zero-shot inference models from the Huggingface website⁶. Considering the NLI is not a binary classification task, we neglect the score of *neutral* prediction and activate the scores of *entailment* and *contradiction* predictions as the final binary output. To help the NLI models better understand the objective of detecting the conspiracy from short texts, we concatenate the input text with the assumption statement (i.e., *That is a conspiracy.*). The model would then give out the answer about whether the assumption statement entails or contradicts the given text. The contradicting answer means that the model predicts the given text as non-conspiracy. The zero-shot test results are in Table 6. The best positive f1-score of 0.57 still does not outperform our proposed conspiracy detection method. Yet, we demonstrate the possibility of those NLI models for the conspiracy detection task.

⁵opt-125m for auto-regressive model and bart-base for the sequence to sequence model

⁶bart-large-mnli, distilbart-mnli-12-1, xlm-roberta-large-xnli

Table 5: Conspiracy detection result of baselines and RoBERTa-based model (R) with different features

Model	accuracy	recall	precision	F1 weighted	F1 negative	F1 positive
Dummy Classifier	0.8000	0.0000	0.0000	0.7111	0.8889	0.0000
Naïve Bayes	0.6825	0.8125	0.3672	0.7141	0.7661	0.5058
Logistics	0.7750	0.7750	0.4627	0.7930	0.8464	0.5794
SVM Linear	0.7675	0.7625	0.4519	0.7863	0.8410	0.5674
combined (R)	0.8575	0.7500	0.6186	0.8624	0.9085	0.6780
transcript (R)	0.7875	0.8250	0.4818	0.8050	0.8542	0.6083
video description (R)	0.8075	0.7000	0.5138	0.8177	0.8740	0.5926
Title (R)	0.8025	0.6875	0.5046	0.8130	0.8707	0.5820
Tags (R)	0.8100	0.6875	0.5189	0.8193	0.8762	0.5914

Table 6: Zero-shot and Few-shot F1-scores for NLI Models. The W stands for weighted f1-score; the N stands for Negative score, and the P stands for Positive score. For the settings column, the ZS stands for zero-shot, while FS-16 stands for Few-shot fine-tuned by 16 data instances.

Model	F1-W	F1-N	F1-P	Setting
bart-large-mnli	0.83	0.91	0.53	ZS
	0.83	0.91	0.54	FS-16
	0.72	0.77	0.53	FS-32
	0.82	0.88	0.57	FS-64
	0.79	0.85	0.54	FS-128
distilbart-mnli-12-1	0.83	0.91	0.52	ZS
	0.83	0.91	0.52	FS-16
	0.65	0.69	0.48	FS-32
	0.77	0.83	0.52	FS-64
	0.83	0.90	0.54	FS-128
xlm-roberta-large-xnli	0.34	0.34	0.34	ZS
	0.32	0.31	0.33	FS-16
	0.40	0.42	0.32	FS-32
	0.65	0.72	0.36	FS-64
	0.62	0.68	0.37	FS-128

Topic Classification

We perform topic inference to understand the type of conspiracy theory of a video published on YouTube. For the ground-truth dataset, we randomly sample 100 videos and label them by the first author based on Table 3.

In Figure 2, we investigate how sensitive our topic inference method is. The method has two parameters: dominance and the minimum number of words matched. Dominance (Zumpe and Michael 1986) is a metric that is commonly used to study the diversity of a community. A higher dominance score suggests a higher percentage of the words matched with one topic (i.e., if dominance is 1, all words matched are in one topic). Hence, having a threshold for dominance to be higher will ensure that the matched topic will have higher accuracy, but lesser videos are likely to be matched. The number of words matched also interplays with the retrieval and accuracy. Figure 2 shows this relationship. For example, when we consider a match of the topic

to be that it requires at least one word matched, 76 videos are matched with a topic, but the accuracy of matching is 0.789. If we increase the threshold to at least 10 words and the dominance score to be greater than 0.6, only 14 videos are matched with a topic but all matching is correct. When there are at least 2 words matched and a dominance threshold of greater than 0.5, the accuracy is 0.842, and 57 videos are matched with a topic.

We explore the topics covered by conspiracy videos using the method outlined above. By using the parameters of at least 2 words matched and a dominance threshold of greater than 0.5, we apply the matching to all the videos with conspiracy theories in our dataset to understand the distribution of the conspiracy topics. Out of the 1,144 conspiracy videos in the dataset, 770 videos have been matched with a topic. Figure 3 shows the distribution of the detected topics. Topics are relatively well distributed, and the top four topics are “Ethnicity, Race, and Religion,” “Government, Politics, and Conflict,” “Science and technology,” and “Extraterrestrials and UFOs.”

Discussion

In this paper, we propose a new dataset, YOUNICON, for the detection of conspiracy theories on YouTube over various topics. While conspiracy theories have been studied for decades across different disciplines, a large-scale dataset of videos on popular social media services will accelerate research on the production and consumption of conspiracy theories on online platforms.

YOUNICON offers a plethora of opportunities to study the subject of conspiracy theories from the text data. First, we hope that the automatic detection of conspiracy theories can be deeply explored by the machine learning community and potentially result in real-world tools to assist and facilitate the work of fact-checkers (e.g., pointing out not only conspiracy theories videos but the exact time the conspiracy theory appears within the video). Our study gives a first step in this direction by exploring standard classification techniques, providing the first assessment of the potential of automated detection of conspiracy theories, and also a baseline for future comparisons. Second, we hope researchers can use the dataset to study the dynamics of conspiracy theories on systems like YouTube. As this dataset contains all videos

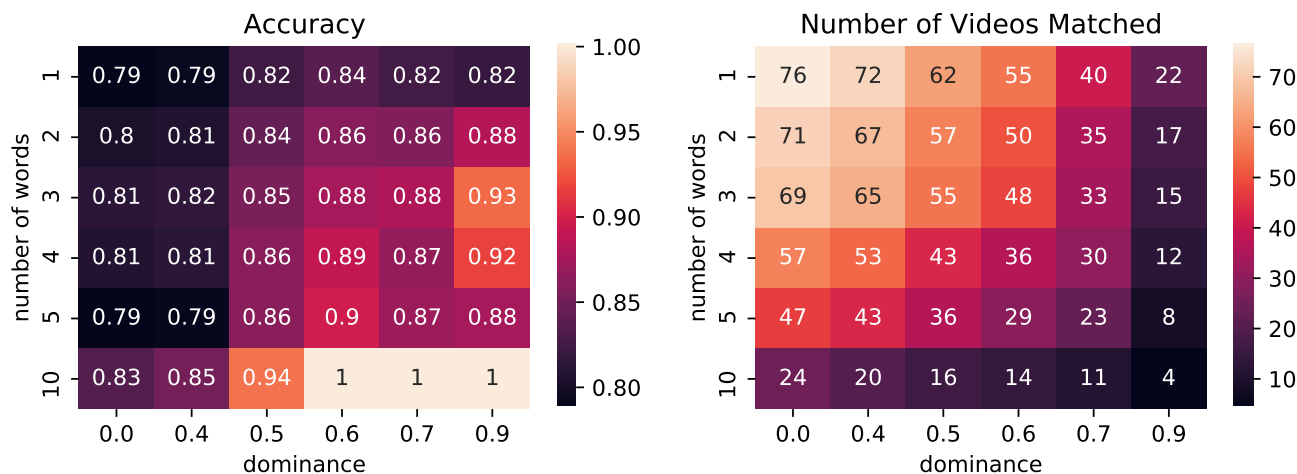


Figure 2: Sensitivity of Number of words and Dominance on Accuracy and Videos Matched

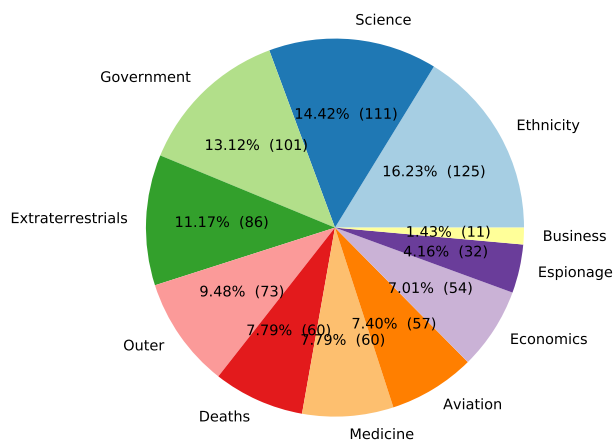


Figure 3: Topic distribution

that are available in the channel’s lifetime (as long as it is not removed from the platform), we are able to study how these content creators have evolved their production strategies over time. For example, do channels focus on a particular type of conspiracy over time or do they adopt a more generalist approach and produce a variety of content? Are there relationships between engagement and topics of conspiracy? The dataset has the potential to answer such questions.

Similarly, for the content consumers (or the video audience), the comments included in the dataset can act as a peephole for us to analyze the behaviour of their consumption of conspiracy theories. In other words, researchers can potentially trace a conspiracy pathway, and look at how people can get involved in the echo chambers of conspiracy theories.

Future works can include looking for better ways to perform topic classification. While Wikipedia’s list of conspiracy theories is used here, this classification can serve as a

starting point for a better taxonomy to be developed. Given the advances of large language models (LLMs), it would be worth exploring the prompting approach with recent LLMs or the in-context learning approach with prompt tuning for the conspiracy detection task.

FAIR Consideration

The proposed dataset follows the FAIR principles of Findability, Accessibility, Interoperability, and Reuse-ability. The dataset can be found and accessed through Zenodo at the DOI: <https://doi.org/10.5281/zenodo.7466262>. Keywords for the topics of conspiracy theories are also shared as a CSV file for the use of other researchers for works related to conspiracy theories. Hence, the data satisfies reusability and interoperability.

Ethical Consideration

We carefully designed our dataset from the data collection period. We collect only publicly available data on YouTube with the use of YouTube’s Data API. Also, our approach is approved by the Institutional Review Board of Singapore Management University (IRB-22-129-A071(922)). To safeguard the interests of our labelers on Amazon Mechanical Turks, they are informed that the content that the conspiracy theories are not true and that withdrawal from the study is without penalty. Helplines are also provided to the participants in the event of any negative emotions.

References

An, J.; Kwak, H.; Lee, C. S.; Jun, B.; and Ahn, Y.-Y. 2021. Predicting Anti-Asian Hateful Users on Twitter during COVID-19. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, 4655–4666.

Bessi, A.; Zollo, F.; Vicario, M. D.; Puliga, M.; Scala, A.; Caldarelli, G.; Uzzi, B.; and Quattrociocchi, W. 2016. Users Polarization on Facebook and Youtube. *PLOS ONE*, 11: e0159641.

- Bowman, S. R.; Angeli, G.; Potts, C.; and Manning, C. D. 2015. A large annotated corpus for learning natural language inference. In *Empirical Methods in Natural Language Processing*.
- Butter, M.; and Knight, P. 2018. The History of Conspiracy Theory Research: A Review and Commentary. In *Conspiracy Theories and the People Who Believe Them*. Oxford University Press. ISBN 9780190844073.
- Center, P. R. 2012. YouTube & News — Pew Research Center.
- Commission, R.; et al. 2020. Royal Commission of Inquiry into the terrorist attack on Christchurch masjidain on 15 March 2019.
- Coninck, D. D.; Frissen, T.; Matthijs, K.; d’Haenens, L.; Lits, G.; Champagne-Poirier, O.; Carignan, M. E.; David, M. D.; Pignard-Cheynel, N.; Salerno, S.; and Généreux, M. 2021. Beliefs in Conspiracy Theories and Misinformation About COVID-19: Comparative Perspectives on the Role of Anxiety, Depression and Exposure to and Trust in Information Sources. *Frontiers in Psychology*, 12.
- Dagnall, N.; Drinkwater, K.; Parker, A.; Denovan, A.; and Parton, M. 2015. Conspiracy theory and cognitive style: A worldview. *Frontiers in psychology*, 6: 206.
- Douglas, K. M.; Sutton, R. M.; and Cichocka, A. 2017. The psychology of conspiracy theories. *Current directions in psychological science*, 26(6): 538–542.
- Enders, A. M.; Uscinski, J. E.; Seelig, M. I.; Klofstad, C. A.; Wuchty, S.; Funchion, J. R.; Murthi, M. N.; Premaratne, K.; and Stoler, J. 2021. The Relationship Between Social Media Use and Beliefs in Conspiracy Theories and Misinformation. *Political Behavior*, 1–24.
- Faddoul, M.; Chaslot, G.; and Farid, H. 2020. A longitudinal analysis of YouTube’s promotion of conspiracy videos. *arXiv preprint arXiv:2003.03318*.
- Fleiss, J. L. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5): 378.
- Galende, B. A.; Hernández-Peñaloza, G.; Uribe, S.; and García, F. Á. 2022. Conspiracy or not? A deep learning approach to spot it on Twitter. *IEEE Access*, 10: 38370–38378.
- Ginossar, T.; Cruickshank, I. J.; Zheleva, E.; Sulskis, J.; and Berger-Wolf, T. 2022. Cross-platform spread: vaccine-related content, sources, and conspiracy theories in YouTube videos shared in early Twitter COVID-19 conversations. *Human vaccines & immunotherapeutics*, 18(1): 1–13.
- Goertzel, T. 1994. Belief in Conspiracy Theories. *Political Psychology*, 15(4): 731–742.
- Goodrow, C. 2017. You know what’s cool? A billion hours.
- Hussein, E.; Juneja, P.; and Mitra, T. 2020. Measuring Misinformation in Video Search Platforms: An Audit Study on YouTube. *ACM Conference on Human-Computer Interaction*, 4(CSCW1).
- Jackson, L. 2017. *Top 201 Conspiracy Theory Videos on YouTube: Full Color Version*. ISBN 9780692995235.
- Jolley, D.; Douglas, K. M.; Leite, A. C.; and Schrader, T. 2019. Belief in conspiracy theories and intentions to engage in everyday crime. *British Journal of Social Psychology*, 58: 534–549.
- Joulin, A.; Grave, E.; Bojanowski, P.; and Mikolov, T. 2016. Bag of Tricks for Efficient Text Classification. *arXiv preprint arXiv:1607.01759*.
- Keeley, B. L. 1999. Of conspiracy theories. *The journal of Philosophy*, 96(3): 109–126.
- Kumar, S.; Asthana, R.; Upadhyay, S.; Upreti, N.; and Akbar, M. 2020. Fake news detection using deep learning models: A novel approach. *Transactions on Emerging Telecommunications Technologies*, 31(2): e3767.
- Kwak, H.; An, J.; and Ahn, Y.-Y. 2020. A systematic media frame analysis of 1.5 million new york times articles from 2000 to 2017. In *12th ACM Conference on Web Science*, 305–314.
- Landis, J. R.; and Koch, G. G. 1977. The measurement of observer agreement for categorical data. *biometrics*, 159–174.
- Ledwich, M.; and Zaitsev, A. 2020. Algorithmic extremism: Examining YouTube’s rabbit hole of radicalization. *First Monday*.
- Lewis, M.; Liu, Y.; Goyal, N.; Ghazvininejad, M.; Mohamed, A.; Levy, O.; Stoyanov, V.; and Zettlemoyer, L. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Association for Computational Linguistics*.
- Lin, X.; Liao, X.; Xu, T.; Pian, W.; and Wong, K.-F. 2019. Rumor Detection with Hierarchical Recurrent Convolutional Neural Network. In *Natural Language Processing and Chinese Computing*.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Mahl, D.; Schäfer, M. S.; and Zeng, J. 2022. Conspiracy theories in online environments: An interdisciplinary literature review and agenda for future research. *new media & society*, 14614448221075759.
- Mahl, D.; Zeng, J.; and Schäfer, M. S. 2021. From “Nasa Lies” to “Reptilian Eyes”: Mapping communication about 10 conspiracy theories, their communities, and main propagators on Twitter. *Social Media+ Society*, 7(2): 20563051211017482.
- Mann, H. B.; and Whitney, D. R. 1947. On a test of whether one of two random variables is stochastically larger than the other. *The annals of mathematical statistics*, 50–60.
- Marcellino, W.; Helmus, T. C.; Kerrigan, J.; Reininger, H.; Karimov, R. I.; and Lawrence, R. A. 2021. *Detecting Conspiracy Theories on Social Media: Improving Machine Learning to Detect and Understand Online Conspiracy Theories*. Santa Monica, CA: RAND Corporation.
- Memon, S. A.; and Carley, K. M. 2020. Characterizing covid-19 misinformation communities using a novel twitter dataset. *arXiv preprint arXiv:2008.00791*.
- Michel, J.-B.; Shen, Y. K.; Aiden, A. P.; Veres, A.; Gray, M. K.; Team, G. B.; Pickett, J. P.; Hoiberg, D.; Clancy, D.;

- Norvig, P.; et al. 2011. Quantitative analysis of culture using millions of digitized books. *science*, 331(6014): 176–182.
- Moffitt, J.; King, C.; and Carley, K. M. 2021. Hunting conspiracy theories during the COVID-19 pandemic. *Social Media+ Society*, 7(3): 20563051211043212.
- Monaci, S. 2021. The pandemic of conspiracies in the covid-19 age: How twitter reinforces online infodemic. *Online Journal of Communication and Media Technologies*, 11.
- Monroe, B. L.; Colaresi, M. P.; and Quinn, K. M. 2008. Fightin' words: Lexical feature selection and evaluation for identifying the content of political conflict. *Political Analysis*, 16(4): 372–403.
- Phillips, S. C.; Ng, L. H. X.; and Carley, K. M. 2022. Hoaxes and Hidden Agendas: A Twitter Conspiracy Theory Dataset: Data Paper. In *Companion Proceedings of the ACM Web Conference*.
- Pummerer, L.; Böhm, R.; Lilleholt, L.; Winter, K.; Zettler, I.; and Sassenberg, K. 2022. Conspiracy theories and their societal effects during the COVID-19 pandemic. *Social Psychological and Personality Science*, 13(1): 49–59.
- Statista. 2022. Most popular social networks worldwide as of January 2022, ranked by number of monthly active users.
- Sutton, R. M.; and Douglas, K. M. 2020. Conspiracy theories and the conspiracy mindset: implications for political ideology. *Current Opinion in Behavioral Sciences*, 34: 118–122.
- Tomlein, M.; Pecher, B.; Simko, J.; Srba, I.; Moro, R.; Stefancova, E.; Kompan, M.; Hrcakova, A.; Podrouzek, J.; and Bielikova, M. 2021. An audit of misinformation filter bubbles on YouTube: Bubble bursting and recent behavior changes. In *ACM Recommender Systems*.
- Uscinski, J. E. 2020. *Conspiracy theories: A primer*. Rowman & Littlefield Publishers.
- Vincent, J. 2020. Something in the air: Conspiracy theorists say 5G causes novel coronavirus, so now they're harassing and attacking UK telecoms engineers.
- Vosoughi, S.; Roy, D.; and Aral, S. 2018. The spread of true and false news online. *science*, 359(6380): 1146–1151.
- Wikipedia contributors. 2022. List of conspiracy theories — Wikipedia, The Free Encyclopedia. [Online; accessed 18-May-2022].
- Williams, A.; Nangia, N.; and Bowman, S. 2018. A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Wood, M. J.; Douglas, K. M.; and Sutton, R. M. 2012. Dead and alive: Beliefs in contradictory conspiracy theories. *Social psychological and personality science*, 3(6): 767–773.
- YouTube. 2022. Understand audience engagement - Computer - YouTube Help.
- Zhang, S.; Roller, S.; Goyal, N.; Artetxe, M.; Chen, M.; Chen, S.; Dewan, C.; Diab, M.; Li, X.; Lin, X. V.; et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.
- Zumpe, D.; and Michael, R. P. 1986. Dominance index: a simple measure of relative dominance status in primates. *American Journal of Primatology*, 10(4): 291–300.