# Black-box Audit of YouTube's Video Recommendation: Investigation of Misinformation Filter Bubble Dynamics

**Matus Tomlein**[1] , **Branislav Pecher**[1,2] , **Jakub Simko**[1] , **Ivan Srba**[1] , **Robert Moro**[1] , **Elena Stefancova**[1] , **Michal Kompan**[1,3] , **Andrea Hrckova**[1] , **Juraj Podrouzek**[1] and **Maria Bielikova**[1,3]

[1]Kempelen Institute of Intelligent Technologies
[2]Brno University of Technology
[3]slovak.AI
[name].[surname]@kinit.com

## Abstract

We investigated the creation and bursting dynamics of misinformation filter bubbles on YouTube using a black-box sockpuppeting audit technique. In this study, pre-programmed agents acting as YouTube users stimulated YouTube's recommender systems: they first watched a series of misinformation promoting videos (bubble creation) and then a series of misinformation debunking videos (bubble bursting). Meanwhile, agents recorded videos recommended to them by YouTube. After manually annotating these recommendations, we were able to quantify the portion of misinformative videos among them. The results confirm the creation of filter bubbles (albeit not in all situations) and show that these bubbles can be bursted by watching credible content. Drawing a direct comparison with a previous study, we do not see improvements in overall quantities of misinformation recommended.

## 1 Introduction

In this work, we investigate the *misinformation filter bubble* creation and bursting on YouTube. In an *auditing study*[1] we simulate user behavior on the YouTube platform, record platform responses (e.g., search results, recommendations) and manually annotate them for the presence of misinformative content. Then, we quantify the dynamics of misinformation filter bubble creation and also dynamics of bubble bursting, which is the novel aspect of the study. Our study adds to the previous works that used *audits* to quantify the misinformative content recommendations and filter bubbles in social media [Hussein *et al.*, 2020; Papadamou *et al.*, 2020].

The general motivation of our work is to emphasize the *need for independent oversight of personalization behavior of large platforms*. In the past, platforms have been accused of being contributors to the misinformation spreading due to their personalization routines. Simultaneously, they have been reluctant to revise these routines [Zuboff, 2019; Vaidhyanathan, 2018]. And when they promise some changes,

there is a lack of effective public oversight that could quantitatively evaluate their fulfillment. Auditing studies are tools that may improve such oversight.

While previous works investigated how a user can enter a filter bubble, no audits have covered *if*, *how* or with what *"effort"* can the user "burst" (exit or lessen) the bubble. Multiple studies demonstrated that watching a series of misinformative videos would strengthen the further presence of such content in recommendations [Abul-Fottouh *et al.*, 2020; Spinelli and Crovella, 2020]. However, no studies investigated what type of user's watching behavior (e.g., switching to credible news videos or conspiracy debunking videos) would be needed to lessen the amount of misinformative content recommended to the user. Such studies would shed more light at the inner workings of YouTube's personalization, but also help improve the social, educational, or psychological resilience strategies against misinformation.

**The first contribution** of this work is the investigation YouTube's personalization behavior in a situation when a user with misinformation promoting watch history (i.e., with a developed misinformation filter bubble) starts to watch content debunking the misinformation (in an attempt to burst that misinformation filter bubble). The key finding is that watching credible content generally improves the situation, albeit with varying effects and forms, mainly depending on particular misinformation topic.

We aligned our methodology with previous works, most notably with the work of Hussein et al. [Hussein *et al.*, 2020] who also investigated the creation of misinformation filter bubbles using user simulation. *As part of our study, we replicated parts of Hussein's study*. We re-used maximum of Hussein's seed data (topics, queries, videos), used similar scenarios and the same data annotation scheme. Therefore, we were able to directly compare the outcomes of both studies. Due to recent changes in YouTube policies [YouTube, 2020], we expected to see less filter bubble creation behavior than Hussein et al. However, this was generally not the case.

**As the second contribution**, we report changes in misinformation video occurrences on YouTube, which took place since the study of Hussein et al. [Hussein *et al.*, 2020] (mid 2019). We observe worse situation regarding the topics of vaccination and (partially) 9/11 conspiracies and some improvements (less misinformation) for moon landing or chemtrails conspiracies.

---

[1]The implementation of the experimental infrastructure and data collected are available at https://github.com/kinit-sk/yaudit-recsys-2021

## 2 Background and related work

*Misinformation filter bubbles* can be defined as states of intellectual isolation in false beliefs or a manipulated perceptions of reality. They can be characterized by a high homogeneity of recommendations/search results that share the same positive stance towards misinformation. The existing studies confirmed the effects of filter bubbles in YouTube recommendations and search results. Spinelli et al. [2020] found that chains of recommendations lead away from reliable sources and towards extreme and unscientific viewpoints. Similarly, Ribeiro et al. [2020] concluded that YouTube's recommendation contributes to further radicalization of users. Abul-Fottouh et al. [2020] confirmed a homophily effect in which anti-vaccine videos were more likely to recommend other anti-vaccine videos than pro-vaccine ones and vice versa.

An *algorithmic audit* is a systematic statistical probing of an online platform, used for quantification of this proportion [Sandvig *et al.*, 2014; Hussein *et al.*, 2020].

*Crowdsourcing audit studies* are conducted using real user data. Silva et al. [2020] developed a browser extension to collect personalized ads with real users on Facebook. Hannak et al. [2013] recruited Mechanical Turk users to run search queries and collected their personalized results. While crowdsourcing audits cover more realistic user conditions, this also means they are noisy (e.g. user behavior is influenced by confirmation bias). Moreover, uncontrolled environment makes comparisons difficult or unfeasible; it is difficult to keep users active; audits also raise several privacy issues.

*Sockpuppeting audits* solve these problems by employing non-human bots that impersonate the behavior of users in a predefined controlled way [Sandvig *et al.*, 2014]. They, however, have their own methodological challenges [Hussein *et al.*, 2020]. First is the selection of appropriate seed data (e.g., the initial activity of bots, search queries). Second, the experimental setup must measure the real influence of the investigated phenomena. At the same time, it must minimize confounding factors and noise (e.g., of name, gender or geolocation). Another challenge is how to appropriately label the presence of the audited phenomena (expert-based/crowdsourced [Hussein *et al.*, 2020; Silva *et al.*, 2020] or automatic labeling [Papadamou *et al.*, 2020]).

Audits can be further distinguished by the social media they are applied on (e.g., social networking sites [Silva *et al.*, 2020; Papadamou *et al.*, 2020; Hussein *et al.*, 2020], search engines [Metaxa *et al.*, 2019; Le *et al.*, 2019; Robertson *et al.*, 2018], e-commerce sites [Juneja and Mitra, 2021]), by adaptive systems being investigated (e.g., recommendations [Hussein *et al.*, 2020; Spinelli and Crovella, 2020; Papadamou *et al.*, 2020], up-next recommendation [Hussein *et al.*, 2020], search results [Papadamou *et al.*, 2020; Hussein *et al.*, 2020; Le *et al.*, 2019; Metaxa *et al.*, 2019; Robertson *et al.*, 2018], autocomplete [Robertson *et al.*, 2018]) and by phenomena being studied (e.g., misinformation [Hussein *et al.*, 2020; Papadamou *et al.*, 2020], political bias [Le *et al.*, 2019; Metaxa *et al.*, 2019], political ads [Silva *et al.*, 2020]). Recently, audits also focused on creation of misinformation filter bubbles [Hussein *et al.*, 2020; Papadamou *et al.*, 2020].

## 3 Study design and methodology

To investigate the dynamics of bursting out of a misinformation filter bubble, we conducted an agent-based sockpuppeting audit study. The study took place on YouTube, but its methodology and implementation can be generalized to any adaptive service, where recommendations can be user-observed.

In the study, we let a series of agents (bots) pose as YouTube users. The agents performed pre-defined sequences of video watches and query searches. They also recorded items they saw: recommended videos and search results. The pre-defined actions were designed to first *invoke the misinformation filter bubble effect* by purposefully watching videos with (or leaning towards) misinformative content. Then, agents tried to *mitigate the bubble effect* by watching videos with trustworthy (misinformation debunking) content. Between their actions, the agents were idle for some time to prevent possible carry-over effects. The degree of how deep inside a bubble the agent is was observed through the number and rank of misinformative videos offered to them.

The secondary outcome is the partial replication of a previous study done by Hussein et al. [Hussein *et al.*, 2020] (denoted onwards as the *reference study*). This replication allowed us to draw direct comparisons between quantities of misinformative content that agents encountered now (March 2021) and during the reference study done in mid 2019.

### 3.1 Research Questions, Hypotheses and Metrics

**RQ1 (comparison to the reference study):** *Has YouTube's personalization behavior changed with regards to misinformative videos since the reference study?* In particular, we seek to validate the following hypothesis:

- **H1.1:** Compared on *SERP-MS* and *normalized score* metrics (see below), we would see better scores (after constructing a promoting watch history) than in the reference study in both search and recommendations (given YouTube's pledges [YouTube, 2020]).

**RQ2 (bubble bursting dynamics):** *How does the effect of misinformation filter bubbles change, when debunking videos are watched?* The "means of bubble bursting" would be implicit user feedback – watching misinformation debunking videos. In particular, we seek to validate the following hypotheses:

- **H2.0:** Watching videos belonging to promoting misinformation stance leads to their increased presence in both search results and recommendations (worse SERP-MS and normalized score metrics).

- **H2.1:** Watching the sequence of misinformation debunking videos after the sequence of misinformation promoting videos will improve the metrics *in comparison to the end of the promoting sequence*.

- **H2.2:** Watching the sequence of misinformation debunking videos after the sequence of misinformation promoting videos will improve the metrics *in comparison to the start of the experiment*.

The metrics we use – *SERP-MS* and *normalized score* – are drawn directly from the reference study. Both metrics

quantify misinformation prevalence in a given list of items (videos), which are annotated as either *promoting* (value 1), *debunking* (value -1) or *neutral* (value 0). The output of both metrics is, similarly, from the $\langle -1, 1 \rangle$ interval. Lists populated mostly with debunking content would receive values close to -1, with promoting close to 1 and with balanced or mostly neutral, close to 0. In other words, a score closer to -1 means better score.

**Normalized score.** A metric computed as average of individual annotations of items present in the list. It is suited for unordered, shorter lists (in our case, recommendations).

**SERP-MS (Search result page misinformation score).** A metric capturing amount of misinformation and its rank. It is suited for longer, ordered lists (in our case, search results). It is computed as $SERP\text{-}MS = \frac{\sum_{r=1}^{n}(x_i*(n-r+1))}{\frac{n*(n+1)}{2}}$, where $x_i$ is annotation value, $r$ search result rank and $n$ number of search results in the list [Hussein *et al.*, 2020].

## 3.2 Experiments scenarios

We let agents interact with YouTube following a *scenario* composed of four phases, as depicted in Figure 1.

*Phase 0: Agent initialization.* At the start of a run, the agent fetches its desired configuration, including the YouTube user account and various controlled variables (the variable values are explained further below). Also, the agent fetches $\tau \in T$, a topic with which it will work (e.g., "9/11"). The agent fetches $V_{prom}$ and $V_{deb}$, which are lists of $n_{prom} = 40$ and $n_{deb} = 40$ most popular videos promoting, respectively debunking, misinformation within topic $\tau$. Afterward, it fetches $Q$, a set of $n_q = 5$ search queries related to the particular $\tau$ (e.g., "9/11 conspiracy"). The agent configures and opens a browser in incognito mode, visits YouTube, logs in using the given user account, and accepts cookies. Finally, the agent creates a neutral baseline by visiting the homepage and saving videos, and performing a search phase. In the *search phase*, the agent randomly iterates through search queries in $Q$, executes each query on YouTube, and saves the search results. To prevent any carry-over effect between search queries, the agent waits for $t_{wait} = 20$ minutes after each query.

*Phase 1 (promoting): Create the filter bubble.* For creating a filter bubble effect, the agent randomly iterates through $V_{prom}$ and "watches" each video for $t_{watch} = 30$ minutes (or less, if the video is shorter). Immediately after watching a video, the agent saves video recommendations on that video's page and visits the YouTube homepage, saving video recommendations listed there as well. After every $f_q = 2$ videos, the agent performs another search phase.

*Phase 2 (debunking): Burst the filter bubble.* The agent follows the same steps as in phase 2. The only difference is the use of $V_{deb}$ instead of $V_{prom}$.

*Phase 3: Tear-down.* In this phase, the agent clears YouTube history (using Google's "my activity" section), making the used user account ready for the next run.

For each selected topic, we run the scenario 10 times (in parallel). This way, we were able to deal with recommendation noise present at the platform. In order to run our experiments multiple times, we used the *reset* (delete all history) button provided by Google instead of creating a new user profile for each run. Before deciding to use the *reset* button in our study, we first performed a short verification study to see whether using this button really deletes the whole history and resets the personalization on YouTube. We randomly selected few topics, from which we manually watched few videos (5 for each). Then, we used the reset button and evaluated the difference between videos appearing on the YouTube homepage, recommendations, and search. We found no carry-over effects.

We needed to set up several attributes of agents (e.g., YouTube user profiles). For *geolocation*, we use N. Virginia to allow for better comparison with the reference study. The date of birth for all accounts was arbitrarily set to 6.6.1990 to represent a person roughly 30 years old. The gender was set as "rather not say" to prevent any personalization based on gender. The names chosen for the accounts were composed randomly of the most common surnames and unisex given names used in the US.

There were also *process parameters* that we needed to keep constant. These include 1) $n_{prom} = 40$ and $n_{deb} = 40$ representing the number of seed videos used in promoting and debunking phases; 2) $t_{watch} = 30$ representing the maximum watching time in minutes for every video; 3) $n_q = 5$ representing the number of queries used; 4) $t_{wait} = 20$ representing the wait time in minutes between query yields and 5) $f_q = 2$ representing the number of videos to watch between search phases.

Values of the *process parameters* greatly influence the total running time and results of the experiment. Yet, determining them was not straightforward given many unknown properties of the environment (first and foremost YouTube's algorithms). For example, prior to the experiment, it was unclear how often we need to probe for changes in recommendations and search result personalizations to answer our research questions.

Therefore, *we run a pre-study in which we determined the best parameter setup*. Measuring the Levenshtein distance between ordered results and overlap of lists of recommended videos we determined to run 10 individual agents for each topic, as we observed instability between repeated runs (e.g., the same configuration yielded $\sim 70\%$ of the same recommended videos). For the $n_{prom}$ and $n_{deb}$ parameters, we observed that in some cases, a filter bubble could be detected after 20 watched videos. Yet in others, it was 30 or more. Due to this inconsistency, we opted to watch 40 videos for a phase. To determine the optimal value of $t_{watch}$, we first calculated the average running time of our seed videos. Most of the videos ($\sim 85\%$) had a running time of about 30 minutes or shorter, so 30 minutes became the baseline value. In addition, we compared the results obtained by watching only 30 minutes with results from watching the whole video regardless of its length, but found no apparent differences.

To determine the number of queries $n_q$ and periodicity of searches $f_q$, we ran the scenario with all seed queries introduced by the reference study and used them after every seed video. We observed that the difference in search results
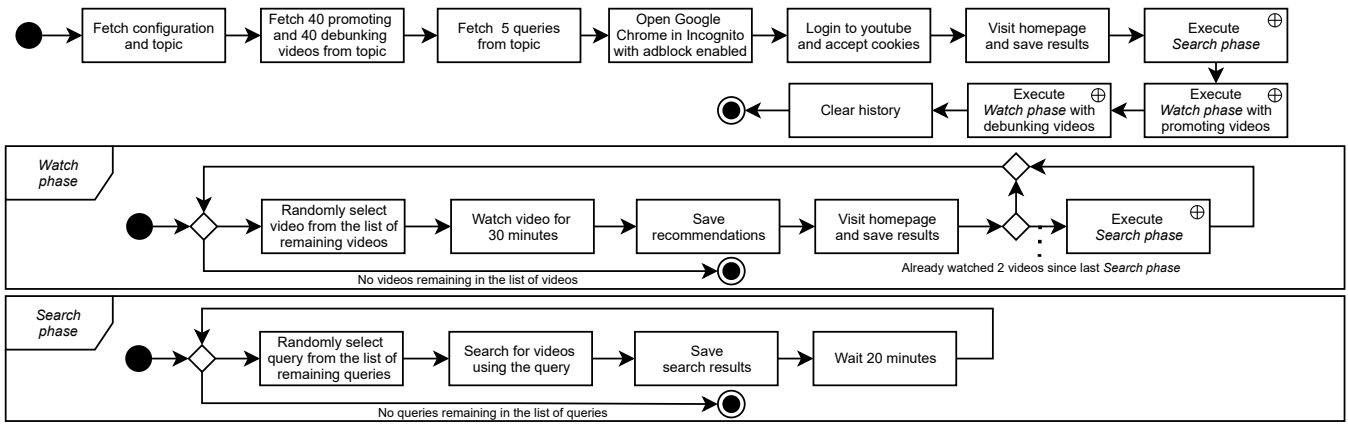
Figure 1: Agent scenario for creating and bursting misinformation filter bubbles

between successive seed videos was not significant. As the choice of search queries and the frequency of their use greatly prolonged the overall running time of the agents, we opted to run the search phase after every second video. In addition, we opted to use only 5 queries per topic.

The only parameter not set by a pre-study is $t_{wait}$, which we set to 20 minutes based on previous studies. These found that the carry-over effect (which we wanted to avoid) is visible for 11 minutes after the search [Hannak *et al.*, 2013; Hussein *et al.*, 2020].

### 3.3 Seed Data

We used 5 topics in our study (same as the reference study): 1) *9/11 conspiracies* claiming that authorities either knew about (or orchestrated) the attack, or that the fall of the twin towers was a result of a controlled demolition, 2)*moon landing conspiracies* claiming the landing was staged by NASA and in reality did not happen, 3) *chemtrails conspiracy* claiming that the trails behind aircraft are purposefully composed of dangerous chemicals, 4) *flat earth conspiracy* claiming that we are being lied to about the spherical nature of Earth and 5) *vaccines conspiracy* claiming that vaccines are harmful, causing various range of diseases, such as autism. The narratives associated with the topics are *popular* (persistently discussed), while at the same time, *demonstrably false*, as determined by the reference study [Hussein *et al.*, 2020].

For each topic, the experiment required two sets of seed videos. The *promoting* set, used to construct a misinformation filter bubble (its videos have a promoting stance towards the conspiratorial narrative or present misinformation). And the *debunking* set, aimed to burst the bubble (and contains videos disproving the conspiratorial narratives).

As a basis for our seed data sets we used data already published in the reference study, which the authors either used as seed data, or collected and annotated. To make sure we use adequate seed data, we re-annotated all of them.

The number of seed videos collected this way was insufficient for some topics (we required twice as many seed videos as the reference study). To collect more, we used an extended version of the seed video identification methodology of the reference study. Following is the list of approaches we used

(in a descending order of priority): YouTube search, other search engines (Google search, Bing video search, Yahoo video search), YouTube channel references, recommendations, YouTube homepage, and known misinformation websites. To minimize any biases, we used a maximum of 3 videos from the same channel.

As for search queries, we required fewer of them than the reference study. We selected a subset based on their popularity on YouTube. Some examples of the used queries are: "9/11 conspiracy", "Chemtrails", "flat earth proof", "anti vaccination", "moon landing fake".

### 3.4 Data collection and annotation

Agents collect videos from three main components on YouTube: 1) *recommendations* appearing next to videos presently watched, 2) *home page* videos and 3) *search results*. In case of recommendations, we collect 20 videos that YouTube normally displays next to a currently watched video (in rare cases, less than 20 videos are recommended). For home page videos and search results, we collect all videos appearing with the given resolution, but no less than 20. In case when less than 20 videos appear, the agent scrolled further down on the page to load more videos.

For each video encountered, the agent collects metadata: 1) *YouTube video ID*, 2) *position* of the video in the list, and 3) *presence of a warning/clarification message* that appears with problematic topics such as COVID-19. Other metadata, such as video *title*, *channel* or *description*, are collected using the YouTube API.

To annotate the collected videos for the presence of misinformation, we used an extended version of the methodology proposed in the reference study. Each video was viewed and annotated by the authors of this study using a code ranging from -1 to 10. The videos are annotated as *debunking* (code -1), when their narrative provides arguments against the misinformation related to the particular topic (such as *"The Side Effects of Vaccines - How High is the Risk?"*), *neutral* (code 0) when the narrative discusses the related misinformation but does not present a stance towards it (such as *"Flat Earthers vs Scientists: Can We Trust Science? — Middle Ground"*), and *promoting* (code 1), when the narrative promotes the related

misinformation (such as *"MIND BLOWING CONSPIRACY THEORIES"*). The codes 2, 3, and 4 have the same meaning as codes -1, 0, and 1, but are used in cases when they discuss misinformation not related to the topic of the run (e.g., video dealing with climate crisis misinformation encountered during a flat earth audit). The code 5 is applied to videos that do not contain any misinformation views (such as *"Gordon's Guide To Bacon"*). This includes completely unrelated videos (e.g., music or reality show videos), but also videos that are related to the general audit topic, but not misinformation (e.g., original news coverage of 9/11 events). In rare cases of videos that are not in English and do not provide English subtitles, code 6 is assigned. Also rare are the cases when the narrative of the video cannot be determined with enough confidence (code 7). Videos removed from YouTube (before they are annotated) are coded as 8. Finally, as an extension of the approach used in the reference study, we use codes 9 and 10 to denote videos that specifically mention misinformation but rather than debunk them, they mock them (9 for related misinformation, 10 for unrelated misinformation, for example *"The Most Deluded Flat Earther in Existence!"*). Mocking videos are a distinct (and often popular) category, which we wanted to investigate separately (however, for the purposes of analysis, they are treated as debunking videos).

To determine how many annotators are needed per video, we first re-annotated the seed videos released by the reference study. Each was annotated by at least two authors, and the annotations were compared between each other and with annotations from the reference study. We achieved Cohen's kappa value of 0.815 between us and 0.688 with the reference study. We identified characteristics of edge cases. Following the re-annotation and the findings from it, when annotating our collected videos, we assign only one annotator per collected video with instructions to indicate and comment if an edge case video is encountered. These were then reviewed by another annotator.

For the purpose of this study and to evaluate our hypotheses, we annotated the following subset of collected videos:

- All recorded *search results*.

- Videos recommended for first 2 seed videos at the start of the run and last 2 seed videos of both phases (resulting in 6 sets of annotated videos per topic). This selection was a compromise between representativeness, correspondence to the reference study, and our capacities.

- We have *not* annotated the *home page videos* for the purpose of this study. These videos were the most numerous, the most heterogeneous, and with little overlap across bots and seed videos.

### 3.5   Data ethics assessment

To consider various ethical issues regarding the research of misinformative content, we carried out a series of data ethics workshops. We explored questions related to data ethics issues [Tranberg *et al.*, 2020] within our audit and its impact on stakeholders. Based on the topics that emerged during the data ethics workshops, we identified different stakeholder groups. The most affected ones were platform users, annotators, content creators, and other researchers. For every stake-

holder group, we devised different engagement strategies and specific action steps. Our main task was to devise countermeasures to the most prominent risks that could emerge for these stakeholder groups.

First, we were concerned about the risk of unjustified flagging of the content as misinformation and their creators as conspirators. To minimize this risk, we decided to report hesitations in the annotation process. These hesitations were consequently back-checked by other annotators and independently validated until the consensus was reached. One of our main concerns was also not to harm or delude other users of the platform. To avoid disproportional boost of the misinformation content by our activity, we select the videos with at least 1000 views and warn annotators not to watch videos online more than one time, or in case of back-checks, two times. After each round, we reset user account and delete the watch history.

Other concerns were connected to the deterioration of well-being of human annotators. Specifically, that their decision-making abilities would be negatively affected after a long annotation process. We proposed the daily routines for annotation, including the breaks during the process and advised to monitor any changes in annotators beliefs. Our annotators also underwent the survey on their tendency to believe in conspiracy theories[2] and none of them showed such tendency at the end of the study.

### 3.6   A note on comparability with the reference study by Hussein et al.

In order to be able to draw comparisons, we kept the methodology of our study as compatible as possible with the previous study by Hussein et al. [Hussein *et al.*, 2020]. We shared the general approach of prompting YouTube with implicit feedback: both studies used similar scenarios of watching a series of misinformation promoting videos and recording search results and recommended videos. We re-used the topics, a subset (for scaling reasons) of search queries, and all available seed videos (complementing the rest by using a similar approach as the reference study). Moreover, both studies used the same coding scheme, metrics, sleep times, and annotated a similar number of videos.

We should also note differences between the studies, which mainly source from different original motivations for our study. For instance, no significant effects of demographics and geolocation of the agents were found in the reference study, so we only controlled these. In Hussein's experiments, all videos were first "watched" and only then all search queries were fired. In our study, we fired all queries after watching *every 2nd* video (with the motivation to get data from the entire run, not just the start and end moment). The reference study created genuine 150 accounts on YouTube, while we used fewer accounts and took advantage of the browsing history reset option. In some aspects, our study had a larger scale: we executed 10 runs for each topic instead of one (to reduce possible noise) and used twice as many seed videos (to make sure that filter bubbles develop). There were also technical differences between the setups, as we used our

---

[2]https://openpsychometrics.org/tests/GCBS/

own implementation of agents (e.g., different browser, ad-blocking software).

Given the methodological alignment (and despite the differences), we are confident to directly compare some of the outcomes of both studies, namely quantity of misinformative content appearing at the end of the promoting phases.

# 4 Results and findings

Following the study design, we executed the study between March 2nd and March 31st, 2021. Together, we executed 50 bot runs (10 for each topic). On average, runs for a single topic took 5 days (bots for a topic ran in parallel). The bots watched 3951 videos (collected 78763 recommendations associated with them, 8526 of them unique), executed 10075 queries (collected 201404 search results, 942 of them unique), and visited homepage 3990 times (collected 116479 videos there, 9977 of them unique). Overall, we recorded 17405 unique videos originating from 6342 channels.

Using the selection strategy and annotation scheme described in Section 3.4, 5 annotators annotated unique 2914 videos (covering 255844 appearances). In total, 244 videos were identified as promoting misinformation (related or unrelated to respective topics), 628 as debunking (including mocking videos), 184 as neutral, 1829 as not about misinformation. Other videos (unknown, non-English, or removed) numbered 29.

We report the results according to research questions and hypotheses defined in Section 3.1. SERP-MS score metrics are reported for search results and mean normalized scores for recommendations. Since the metrics are not normally distributed with some samples of unequal sizes, we make use of non-parametric statistical tests. Pairwise tests are performed using two-sided Mann-Whitney U test. In cases where multiple comparisons by topics are performed, Bonferroni correction is applied on the significance level (in that case $\alpha = 0.05$ is divided by number of topics $n_T = 5$, resulting in $\alpha = 0.01$).

## 4.1 RQ1: Has YouTube's personalization behavior changed since the reference study?

Overall, we see a small change in the mean SERP-MS score across the same search queries in our and reference data: mean SERP-MS worsened from -0.46 (std 0.42) in reference data to -0.42 mean (std 0.3) in our data. However, the distributions are not statistically significantly different (n.s.d.). There is a similar small change towards the promoting spectrum in up-next (first result in recommendation list) and top-5 recommendations (following 5 recommendations). We compared the up-next and top-5 recommendations together (as top-6 recommendations) using last 10 watched promoting videos in reference watch experiments and last two watched videos in our promoting phase. We see mean normalized score worsened from -0.07 (std 0.27) in reference data to -0.04 (std 0.31) in our data. These distributions are also not significantly different (U=45781.5, n.s.d.).

More considerable shifts in the data can be observed when looking at individual topics. Table 1 shows a comparison of SERP-MS scores for top-10 search results between our

and reference data. Improvement can be seen within certain queries for the chemtrails conspiracy that show a large decrease in the number of promoting videos. The reference study reported that this topic receives significantly more misinformative search results compared to all other topics. In our experiments, their proportion was lower than in the 9/11 conspiracy. On the other hand, search results for flat earth conspiracy worsened. Queries such as "flat earth british" resulted in more promoting videos, likely due to new content on channels with similar names. Within the anti-vaccination topic, there is an increase in neutral videos (from 12% to 35%) and thus a drop in debunking videos (from 85% to 61%). This may relate to new content regarding COVID-19.

Table 2 shows a comparison of normalized scores for up-next and top-5 recommendations. Only the moon landing and anti-vaccination topics come from statistically significantly different distributions. Similar to search results, recommendations for the 9/11 and anti-vaccination conspiracy topics worsened. There were more promoting videos on the 9/11 topic (27% instead of 18%). In the anti-vaccination topic, we observed a drop in debunking videos (from 29% to 9%) and a subsequent increase in neutral (from 70% to 78%) and promoting videos (from 1% to 8%). The change within the anti-vaccination controversy is even more pronounced when looking at up-next recommendations separately. Within up-next, the proportion of debunking videos drops from 77% to 19%, neutral videos increase from 22% to 70%, and promoting increase from 1 to 11%. On the other hand, in the moon landing topic, we see much more debunking video recommendations—40% instead of 23% in reference data.

These results bring up a need to distinguish between *endogenous* (changes in algorithms, policy decisions made by platforms to hide certain content) and *exogenous* factors (changes in content, external events, behavior of content creators) as discussed by Metaxa et al. [Metaxa *et al.*, 2019]. Our observations show that search results and recommendations were in part influenced by exogenous changes in content on YouTube. Within the chemtrails conspiracy, we observed results related to a new song by Lana del Rey that mentions "Chemtrails" in its name. Search results and recommendations in the anti-vaccination topic seem to be influenced by COVID-19. Flat earth conspiracy videos were influenced by an increased amount of activity within a single conspiratorial channel.

## 4.2 RQ2: What is the effect of watching debunking videos after the promoting phase?

Answering this question requires three comparisons:

1. comparison of metrics between start of promoting phase (S1) and end of promoting phase (E1),

2. comparison of metrics between end of promoting phase (E1) and end of debunking phase (E2),

3. comparison of metrics between start of promoting phase (S1) and end of debunking phase (E2).

*Comparison (1)* shows changes in search results and recommendations after watching promoting videos (E1) compared to the start of the experiment (S1). If there was a

Table 1: Comparison of SERP-MS scores for top-10 search results with data from the reference study. The scores range from $\langle -1, 1 \rangle$, where -1 denotes a debunking and 1 a promoting stance towards the conspiracy. Only search results from queries that were executed both by the reference study and us are considered.

| Topic | Hussein | Ours | Change | Inspection |
|---|---|---|---|---|
| 9/11 | -0.16 | -0.06 | No (n.s.d.) | Smaller changes that depend on search query. |
| Chemtrails | -0.2 | -0.47 | No (n.s.d.) | Drop in promoting videos (from 45% to 12%) in 2 queries. |
| Flat earth | -0.58 | -0.41 | No (n.s.d.) | 2 queries worsen a lot due to new content. Other queries improve. |
| Moon landing | -0.6 | -0.59 | No (n.s.d.) | Smaller decrease in number of neutral and increase of debunking videos. |
| Anti-vaccination | -0.8 | -0.63 | Worse (U=324,p=1.3e−9) | Drop in number of debunking and increase in number of neutral videos. |

Table 2: Comparison of normalized scores for up-next and top-5 recommendations with data from the reference study. Normalized scores range from $\langle -1, 1 \rangle$, where -1 denotes a debunking and 1 a promoting stance towards the conspiracy. Last 10 out of 20 watched videos in reference data are considered. Last 2 out of 40 watched videos in our data are considered.

| Topic | Hussein | Ours | Change | Inspection |
|---|---|---|---|---|
| 9/11 | 0.14 | 0.26 | No (n.s.d.) | Similar distribution, more promoting videos. |
| Chemtrails | 0.05 | 0.03 | No (n.s.d.) | More neutral results. |
| Flat earth | -0.16 | -0.15 | No (n.s.d.) | Similar distribution. |
| Moon landing | -0.08 | -0.32 | Better (U=2954.5,p=8e−6) | More debunking videos. |
| Anti-vaccination | -0.28 | -0 | Worse (U=664,p=1.6e−9) | Less debunking videos, more neutral and promoting. |

Table 3: Comparison of SERP-MS scores for top-10 search results in promoting and debunking phase of our experiment. Three points are compared: start of promoting phase (S1), end of promoting phase (E1), end of debunking phase (E2).

| Topic | SERP-MS | Change | | Inspection |
|---|---|---|---|---|
| 9/11 | S1: -0.07 E1: -0.06 E2: -0.11 | S1–E1: n.s.d. E1–E2: n.s.d. S1–E2: n.s.d. | | E2: More debunking videos in one query (30% instead of 12% at S1 and 11% at E1 in query "9/11"). |
| Chemtrails | S1: -0.45 E1: -0.47 E2: -0.49 | S1–E1: n.s.d. E1–E2: n.s.d. S1–E2: (U=915,p=0.0097) | better | E2: The "Chemtrail" search query showed an increase in number of debunking videos (from 66% at S1 and 69% at E1 to 80%) and a decrease in promoting (from 10% to 0%). |
| Flat earth | S1: -0.27 E1: -0.41 E2: -0.45 | S1–E1: (U=762.5,p=0.0004) E1–E2: n.s.d. S1–E2: (U=704.5,p=0.0001) | better better | E1: Change goes against expectations. Promoting videos disappear in 3 search queries and decrease in another one (from 36% to 30%). E2: Similar change as in E1 with a further decrease in promoting videos in one query (from 30% to 22%) and reordered videos in another. |
| Moon landing | S1: -0.57 E1: -0.57 E2: -0.59 | S1–E1: n.s.d. E1–E2: n.s.d. S1–E2: (U=900,p=0.0068) | better | E2: Reordered search results in "moan hoax" query—debunking videos moved higher. |
| Anti-vacc. | S1: -0.6 E1: -0.63 E2: -0.68 | S1–E1: n.s.d. E1–E2: (U=699.5,p=0.0054) S1–E2: (U=641.5,p=0.0001) | better better | E2: Increase in debunking videos across multiple queries (from 60% at S1 and 61% at E1 to 67%). |

Table 4: Comparison of changes in average normalized scores for top-10 recommendations in promoting and debunking phase of our experiment. Three points are compared: start of promoting phase (S1), end of promoting phase (E1), end of debunking phase (E2).

| Topic | Score | Change | | Inspection |
|---|---|---|---|---|
| 9/11 | S1: 0.1<br>E1: 0.42<br>E2: 0.07 | S1–E1:<br>(U=45.5,p=2.6e−5)<br>E1–E2:<br>(U=28,p=2.9e−6)<br>S1–E2: n.s.d. | worse<br><br>better | E1: Number of promoting videos increased (from 14% to 43%) and neutral videos decreased (from 83% to 56%).<br>E2: The numbers of promoting and neutral videos returned to levels comparable to start (13% and 82%). |
| Chemtrails | S1: 0<br>E1: 0.05<br>E2: -0.15 | S1–E1: n.s.d.<br>E1–E2: better (U=323, p=0.0006)<br>S1–E2: better (U=330, p=0.0002) | | E2: There is an increase in a number of debunking videos (from 0% at S1 and 3% at E1 to 19%). In return, we end up in a state that is better than at the start. |
| Flat earth | S1: -0.17<br>E1: -0.06<br>E2: -0.47 | S1–E1: n.s.d.<br>E1–E2: better (U=375, p=1.8e−6)<br>S1–E2: better (U=347, p=0.0001) | | E2: Similar to the Chemtrails conspiracy, there is an increase in number of debunking videos (from 19% at S1 and 16% at E1 to 48%). |
| Moon landing | S1: -0.2<br>E1: -0.4<br>E2: -0.42 | S1–E1: n.s.d.<br>E1–E2: n.s.d.<br>S1–E2: n.s.d. | | E1: Mean normalized scores changes against expectation and improves (but not significantly). |
| Anti-vacc. | S1: -0.1<br>E1: 0.04<br>E2: -0.37 | S1–E1:<br>(U=74.5,p=0.0008)<br>E1–E2:<br>(U=310,p=2.5e−6)<br>S1–E2:<br>(U=307.5,p=0.0002) | worse<br><br>better<br><br>better | E1: Increase in number of promoting videos (from 2% to 13%).<br>E2: Increase of debunking videos (from 12% at S1 and 9% at E1 to 37%) and disappearance of promoting (from 2% at S1 and 13% at E1 to 0%). |

misinformation bubble created, we would expect the metrics to worsen due to watching promoting videos. Regarding search results, the distribution of SERP-MS scores between S1 and E1 is indeed significantly different (MW U=34118.5, p-value=0.028). However, the score actually improves—mean SERP-MS score changed from -0.39 (std 0.28) to -0.42 (std 0.3). Table 3 shows the change for individual topics. Only the flat earth conspiracy shows significant differences and improved the SERP-MS score due to a decrease in promoting and an increase of debunking videos. Top-10 recommendations also change their distribution of normalized scores significantly at E1 compared to S1 (MW U=4085, p-value=0.0397). We observe that the mean normalized score worsens from -0.07 (std 0.24) to 0.01 (std 0.31). Looking at individual topics in Table 4, we can see that the change is significant in topics 9/11 and anti-vaccination that gain more promoting videos.

*Comparison (2)* relates the change in search results and recommendations between the end of promoting phase (E1) and the end of debunking phase (E2). We expect the metrics would improve due to watching debunking videos, i.e., that we would observe misinformation bubble bursting. However, SERP-MS scores in search results between E1 and E2 are not from statistically significantly different distributions, which is consistent with the fact that we did not observe misinformation bubble creation in search results in the first place. Table 3 shows that only a single topic—anti-vaccination—significantly changed its distribution and improved its mean score. Nevertheless, we see minor improvements in SERP-

MS scores also in other topics. Top-10 recommendations show more considerable differences and their overall distribution is significantly different comparing E1 and E2 (MW U=7179.5, p-value=1.8e−9). Mean normalized score improves from 0.01 (std 0.31) to -0.27 (std 0.27). Table 4 shows significantly different distributions for all topics except for moon landing conspiracy. All topics show an improvement in normalized scores. The 9/11 topic shows a decrease in promoting videos, while other topics show an increase in the number of debunking videos.

*Comparison (3)* shows differences between the start (S1) and end of the experiment (E2). We expect the metrics would improve due to watching debunking videos despite watching promoting videos before that. The distribution of SERP-MS scores in search results is statistically significantly different when comparing S1 and E2 (MW U=36515, p-value=0.0002). Overall, we see an improvement in mean SERP-MS score from -0.39 (std 0.28) to -0.46 (std 0.29). In contrast with comparison (2), Table 3 shows that all topics except 9/11 significantly changed their distributions. All topics show an improvement according to our expectations. The improvement is due to increases in debunking videos, decreases in promoting videos, or reordered search results in some search queries. Similarly, top-10 recommendations at E2 come from a significantly different distribution than at S1 (MW U=6940.5, p-value=2.9e−7). Mean normalized score improves from -0.07 (std 0.24) to -0.27 (std 0.27). Table 4 shows a significant difference in distributions for all topics except for 9/11 and moon landing conspiracies. Mean nor-

malized scores improve compared to S1 in all topics except for 9/11. Nevertheless, the numbers of promoting and neutral videos in 9/11 topic at E2 are comparable to S1. Other topics show increases in the numbers of debunking videos.

## 5 Discussion and Conclusions

In the paper, we presented an audit of misinformation present in search results and recommendations on the video-sharing platform YouTube. To support reproducibility, we publish the collected data and source codes for the experiment.

We aimed at verifying a hypothesis that there is less misinformation present in both search results and recommendations after recent changes in YouTube policies [YouTube, 2020] (H1.1). The comparison was done against a study done in mid 2019 by Hussein et al. [Hussein *et al.*, 2020]. We were interested, whether we could still observe the formation of misinformation bubbles after watching videos promoting conspiracy theories (H2.0). In contrast to the previous studies, we also examined bubble bursting behavior. Namely, we aimed to verify whether misinformation bubbles could be burst if we watched videos debunking conspiracy theories (H2.1). We also hypothesized that watching debunking videos (even after a previous sequence of promoting videos) would still decrease the amount of misinformation compared to the initial state with no watch history at the start of the study (H2.2).

Regarding hypothesis H1.1, we did not find a significantly different amount of misinformation in search results in comparison to the reference study. A single topic (anti-vaccination) showed a statistically significant difference. However, it did not agree with the hypothesis as the metric *worsened* due to more neutral and less debunking videos. Recommendations showed significant differences across multiple topics but were not significantly different overall. A single topic (moon landing) improved normalized scores of recommendation in agreement with the hypothesis. Yet, the anti-vaccination topic worsened its scores. We suspect the changes in search results and recommendations were influenced mostly by changes in content. Overall, our results did not show a significant improvement in the fight against misinformation on the platform, as stated in the hypothesis.

We did not observe the creation of misinformation filter bubbles in search results (H2.0) despite watching promoting videos. On the other hand, recommendations behaved according to our hypothesis, and their overall normalized scores worsened. Since there was no filter bubble creation effect in search results, we did not observe any bubble bursting effect there. Results did not show a statistically significant difference between the end of promoting phase and the end of the debunking phase. Only a single topic (anti-vaccination) showed a statistically significant difference and an improvement following the hypothesis H2.1. Recommendations showed more considerable differences that were statistically significant and confirmed the hypothesis. Lastly, we showed that watching debunking videos decreases the number of misinformation videos both in search results and recommendations, which confirms our hypothesis H2.2. We observed an improvement of SERP-MS scores in all topics except for one and an improvement of normalized scores for recommendations in most topics.

Based on our results, we can conclude that users, even with a watch history of promoting conspiracy theories, do not get enclosed in a misinformation filter bubble *when they search* on YouTube. However, we do observe this effect in video recommendations with varying degrees depending on the topic. However, *watching debunking videos helps in practically all cases* to decrease the amount of misinformation that the users see. Additionally, although we expected to see less misinformation than the previous studies reported, this was in general not the case. Worsening in the anti-vaccination topic was partially expected due to the COVID-19 pandemic. However, it is interesting that we also observed a worse situation with the 9/11 topic. In fact, this topic served as a sort of a gateway to misinformation videos on other topics.

A limitation of our results lies with the limited amount of topics that we investigated – these did not include, for example, recent QAnon conspiracy and COVID-19 related conspiracies were present only through anti-vaccination narratives. However, our topics were explicitly selected to allow comparison with the reference study. Next, we included only a limited set of agent interactions with the platform (search and video watching). Real users also like or dislike videos, subscribe to channels, leave comments or click on the search results or recommendations. A more human-like bot simulation, with these interactions and possible inclusion of human biases bursting remains our future work.

Nevertheless, our audit showed that YouTube (similar to other platforms), despite their best efforts so far, can still promote misinformation seeking behavior to some extent. The results also motivate the need for independent continuous and automatic audits of YouTube and other social media platforms [Simko *et al.*, 2021], since we observed that the amount of misinformation in a topic could change over time due to endogenous as well as exogenous factors.

## Acknowledgments

## References

[Abul-Fottouh *et al.*, 2020] Deena Abul-Fottouh, Melodie Yunju Song, and Anatoliy Gruzd. Examining algorithmic biases in youtube's recommendations of vaccine videos. *Int. Journal of Medical Informatics*, 140:104175, 2020.

[Hannak *et al.*, 2013] Aniko Hannak, Piotr Sapiezynski, Arash Molavi Kakhki, Balachander Krishnamurthy, David Lazer, Alan Mislove, and Christo Wilson. Measuring personalization of web search. In *Proc. of the 22nd International Conference on World Wide Web (WWW '13)*, page 527–538, New York, NY, USA, 2013. ACM.

[Hussein *et al.*, 2020] Eslam Hussein, Prerna Juneja, and Tanushree Mitra. Measuring misinformation in video search platforms: An audit study on youtube. *Proc. ACM Hum.-Comput. Interact.*, 4(CSCW1), May 2020.

[Juneja and Mitra, 2021] Prerna Juneja and Tanushree Mitra. Auditing E-Commerce Platforms for Algorithmically Curated Vaccine Misinformation. In *Proc. of the 2021 CHI Conference on Human Factors in Computing Systems (CHI '21)*, 2021.

[Le *et al.*, 2019] Huyen Le, Andrew High, Raven Maragh, Timothy Havens, Brian Ekdale, and Zubair Shafiq. Measuring political personalization of Google news search. In *Proc. of the World Wide Web Conference (WWW '19)*, pages 2957–2963, 2019.

[Metaxa *et al.*, 2019] Danaë Metaxa, Joon Sung Park, James A. Landay, and Jeff Hancock. Search media and elections: A longitudinal investigation of political search results. *Proc. ACM Hum.-Comput. Interact.*, 3(CSCW), November 2019.

[Papadamou *et al.*, 2020] Kostantinos Papadamou, Savvas Zannettou, Jeremy Blackburn, Emiliano De Cristofaro, Gianluca Stringhini, and Michael Sirivianos. "it is just a flu": Assessing the effect of watch history on youtube's pseudoscientific video recommendations, 2020.

[Ribeiro *et al.*, 2020] Manoel Horta Ribeiro, Raphael Ottoni, Robert West, Virgílio A. F. Almeida, and Wagner Meira. Auditing radicalization pathways on youtube. In *Proc. of the 2020 Conference on Fairness, Accountability, and Transparency*, page 131–141, New York, NY, USA, 2020. ACM.

[Robertson *et al.*, 2018] Ronald E. Robertson, David Lazer, and Christo Wilson. Auditing the personalization and composition of politically-related search engine results pages. pages 955–965, 2018.

[Sandvig *et al.*, 2014] Christian Sandvig, Kevin Hamilton, Karrie Karahalios, and Cedric Langbort. Auditing algorithms: Research methods for detecting discrimination on internet platforms. *Data and discrimination: converting critical concerns into productive inquiry*, 22:4349–4357, 2014.

[Silva *et al.*, 2020] Márcio Silva, Lucas Santos de Oliveira, Athanasios Andreou, Pedro Olmo Vaz de Melo, Oana Goga, and Fabricio Benevenuto. Facebook ads monitor: An independent auditing system for political ads on facebook. In *Proc. of The Web Conference (WWW '20)*, page 224–234, New York, NY, USA, 2020. ACM.

[Simko *et al.*, 2021] Jakub Simko, Matus Tomlein, Branislav Pecher, Robert Moro, Ivan Srba, Elena Stefancova, Andrea Hrckova, Michal Kompan, Juraj Podrouzek, and Maria Bielikova. Towards continuous automatic audits of social media adaptive behavior and its role in misinformation spreading. In *Adjunct Proc. of the 29th ACM Conference on User Modeling, Adaptation and Personalization (UMAP '21)*, page 411–414, New York, NY, USA, 2021. ACM.

[Spinelli and Crovella, 2020] Larissa Spinelli and Mark Crovella. How youtube leads privacy-seeking users away from reliable information. In *Adjunct Publication of the 28th ACM Conference on User Modeling, Adaptation and Personalization*, page 244–251, New York, NY, USA, 2020. ACM.

[Tranberg *et al.*, 2020] Pernille Tranberg, Gry Hasselbalch, Catrine S. Byrne, and Birgitte K. Olsen. *DATAETHICS – Principles and Guidelines for Companies, Authorities & Organisations*. Dataethics.eu, 2020.

[Vaidhyanathan, 2018] Siva Vaidhyanathan. *Antisocial media: How Facebook disconnects us and undermines democracy*. Oxford University Press, 2018.

[YouTube, 2020] YouTube. Managing harmful conspiracy theories on youtube, 2020.

[Zuboff, 2019] Shoshana Zuboff. *The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power*. Profile Books, 2019.