

A Dataset to Assess Microsoft Copilot Answers in the Context of Swiss, Bavarian and Hessian Elections

Salvatore Romano^{1, 2}, Riccardo Angius¹, Natalie Kerby¹, Paul Bouchaud^{1, 3, 4},
Jacopo Amidei², Andreas Kaltenbrunner^{2, 5}

¹AI Forensics, Paris, France

²AI and Data for Society, IN3, Universitat Oberta de Catalunya, Barcelona, Spain

³Center for Social Analysis and Mathematics, Paris, France

⁴Complex Systems Institute of Paris, France

⁵ISI Foundation, Turin, Italy

salvatore@aiforensics.org, riccardo@aiforensics.org, natalie@aiforensics.org, paul@aiforensics.org,
jamidei@uoc.edu, akaltenbrunner@uoc.edu

Abstract

This study describes a dataset that allows to assess the emerging challenges posed by Generative Artificial Intelligence when doing Active Retrieval Augmented Generation (RAG), especially when summarizing trustworthy sources on the Internet. As a case study, we focus on Microsoft Copilot, an innovative software that integrates Large Language Models (LLMs) and Search Engines (SE) making advanced AI accessible to the general public. The core contribution of this paper is the presentation of the largest public database to date of RAG responses to user prompts, collected during the 2023 electoral campaigns in Switzerland, Bavaria and Hesse. This dataset was compiled with the assistance of a group of experts who posed realistic voter questions and conducted fact-checking of Microsoft Copilot's responses. It contains prompts and answers in English, German, French and Italian. All the collection happened during the electoral campaign, between 21 August 2023 and 2 October 2023. The paper makes available the full set of 5,561 pairs of prompts and answers, including the URLs referenced in the answers. In addition to the dataset itself, we provide 1374 answers labelled by a group of experts who rated the accuracy of the answers in providing factual information, showing that almost one out of three times the chatbot responded with either factually incorrect information or completely nonsensical answers. This resource is intended to facilitate further research into the performance of LLMs in the context of elections, defined as a "high-risk scenario" by the Digital Services Act (DSA) Article 34(1)(c).

Introduction

Microsoft Copilot (Microsoft 2023b), previously branded as Bing Chat, is a conversational AI tool (or chatbot) released to the general public by Microsoft in February 2023 on its search engine Bing. At the moment this paper is written it is also embedded across Microsoft products such as the Microsoft Office suite and the Windows 11 operating system. This AI tool generates answers based on live retrieved search results by combining an LLM with search engine capabilities.

Copyright © 2024, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

In this study, we detail a curated dataset that we developed to evaluate the accuracy and comprehensiveness of Microsoft Copilot's responses to inquiries pertaining to the Swiss general elections on October 8, 2023, and the Bavarian and Hessian regional state elections in Germany on October 22, 2023. Our methodology involved querying the chatbot with a range of questions focused on candidates, polling data, and voting details, alongside broader queries soliciting voting recommendations for specific policy areas such as environmental concerns. The dataset encompasses the chatbot's responses collected from August 21, 2023, to October 2, 2023, offering insights into its factual reliability and informative value in a politically charged context.

Selecting Microsoft Copilot as the focus for this dataset facilitates a comprehensive analysis to address issues like:

- Assessing the impact of SEs powered by LLMs on the fidelity of information disseminated during electoral campaigns.
- Evaluating Microsoft Copilot's capacity to provide precise, timely, and non-partisan insights regarding candidates and political parties.
- Examining the consistency of the provided information over time, across different languages, and concerning specific local contexts.

Background and Related Work

Microsoft Copilot Overview

The chat functionality in Microsoft Copilot is powered by OpenAI's GPT-4 (Achiam et al. 2023), a leading LLM. Key to its operation is the integration of what Microsoft terms the "Prometheus model" (Microsoft 2023). This innovative model merges the advanced processing abilities of OpenAI's GPT with the sophisticated indexing and ranking algorithms of Bing Search. Rather than merely ranking web pages, Prometheus actively parses and analyzes them. This enables it to extract and produce a summary of the information it classifies as relevant, thereby enriching its responses with contextually significant content drawn directly from the web.

The method of enhancing LLMs by incorporating information from external knowledge sources and improving their accuracy with factual information is generally known as Retrieval Augmented Generation (RAG) (Lewis et al. 2021).

In the case of Microsoft Copilot, the model's functionalities are not transparent, as no public and thorough technical description exists of how the tools work together. While a traditional search engine typically lists sources in response to a query (Microsoft 2023), Microsoft Copilot compiles information from different sources into an answer which is designed to appear based on trusted information. Since Copilot can conduct web searches, it is marketed as a tool for delivering "better search [and] more complete answers" (Mehdi 2023). The chatbot provides users with web sources for specific information, but it is also able to generate creative content such as stories and poems based on prompts. Furthermore, Microsoft Copilot users can choose between various languages.

Microsoft Copilot may provide answers to users' questions that include up-to-date information, a special feature not included in the current free version of OpenAI ChatGPT. As Microsoft does not require a subscription to access Copilot, users are allowed to generate five prompts without being logged in, and more if using an account.

Microsoft Copilot's most important feature is its advanced natural language processing technology. But there are concerns about its potential to perpetuate biases and the generated answers' factual correctness per se (Birhane et al. 2023). This potentially dangerous flaw is often referred to as "hallucinating". LLMs string words together based on probability, not based on factual consistency (Bender et al. 2021). With RAG, the LLMs can incorporate information from external knowledge sources, augmenting the accuracy in the task of answering at factual QA (Kwiatkowski et al. 2019). At the same time, the LLM could introduce mistakes in text summarization, and the best way to prevent this is still under discussion (Jiang et al. 2023). This risk raises additional questions given Microsoft's statements on the production of "more than 1.9 billion Copilot chats" in 2023 (Microsoft 2023c).

DSA and Systemic Risks During Elections

By ranking publicly accessible information, search engines like Google and Bing have gained substantial power. For many people, search engines are their preferred source of information on the Internet. As the public debate increasingly takes place online, this has serious consequences for the integrity of elections (Epstein and Robertson 2015; Epstein and Li 2023).

Lawmakers in the EU have recognized the need for regulation of large platforms and search engines. The EU's Digital Services Act (Regulation (EU) 2022/2065), a law introduced in 2022 to regulate digital platforms, requires "very large online platforms" (VLOPs) and "very large search engines" (VLOSEs) with more than 45 million users within the EU to carry out so-called risk assessments and develop mechanisms to mitigate the risks posed by their services. The European Commission has categorized Google Search

and Microsoft Bing as such VLOSE (European Commission 2023). The law explicitly mentions negative effects on the integrity of electoral processes and public debate, as well as on the spread of misinformation, as "systemic risks" that can emanate from Microsoft Bing, Google, and other search engines in Article 34(C) DSA. The providers must thus examine if their services work properly, and take action otherwise. A "systemic risk" is not clearly defined. Still, under the DSA, VLOSEs are obliged to publish transparency reports regularly. The first Bing transparency report (Microsoft 2023a), published on 6 November 2023, mentions Microsoft Copilot once, without elaborating on Microsoft's strategy to mitigate the risk to the integrity of elections caused by the integration of a generative AI feature in its search engine.

This paper introduces a pioneering database designed to assess the unique risks to elections posed by VLOSE combined with Generative AI, in compliance with the DSA guidelines. Additionally, it serves as the initial database for evaluating the accuracy and reliability of information provided by Microsoft Copilot in electoral scenarios.

Electoral Context

This investigation covers three 2023 elections as case studies: the Swiss Federal election of October 22, as well as the state elections in the German federal states of Hesse and Bavaria of October 8.

These were some of the first elections to take place in Germany and Switzerland after the earliest introduction of Microsoft Copilot as Bing Chat. These case studies enable the analysis of different local contexts and political systems, as well as a comparison across different languages (German and English for Germany, as well as German, French, Italian and English for Switzerland).

Switzerland Switzerland's status as a multilingual direct democracy presents a unique and complex electoral context, making it an intriguing case for examining election integrity risks. The Swiss electorate, characterized by its diversity in language and culture, participates in a robust democratic process where they are frequently called upon to vote at the national, cantonal, and communal levels. This frequent engagement in a wide range of issues, from national referendums to local matters, underscores the critical need for voter access to reliable and comprehensible information.

The multilingual nature of the country, with German, French, Italian, and Romansh as national languages, adds another layer of complexity to the information dissemination and verification. Ensuring that accurate and unbiased information is equally accessible across these linguistic divides is paramount for maintaining election integrity.

Moreover, the decentralized nature of Swiss politics, with significant autonomy granted to cantons, further complicates the electoral landscape. Each canton can have its own distinct political culture and voting procedures, which means that understanding and monitoring election integrity requires a nuanced, region-specific approach. In this context, the role of technology and AI tools like Microsoft Copilot in providing information becomes even more critical, as they must cater to a wide array of linguistic, cultural, and regional

needs while upholding the principles of accuracy and impartiality.

Bavaria and Hesse The electoral dynamics in Bavaria and Hesse, two pivotal states in Germany, offer a diverse and enlightening view of the processes and outcomes in the elections for their regional parliaments.

The 2023 Bavarian state election, held on October 8, 2023, was marked by significant political dynamics. The Christian Social Union (CSU) maintained its dominant position, albeit with a marginal decrease in its vote share, illustrating its continued influence in Bavaria's political landscape, rooted in regional issues and traditional values. The Free Voters of Bavaria (FW) made notable gains, securing second place, a testament to their increasing popularity. The election also highlighted the rise of the Alternative for Germany (AfD), emphasizing a shift in the political spectrum. These results demonstrate the complex interplay of regional identity, traditional conservatism, and emerging political forces in shaping Bavarian politics.

The 2023 Hessian state election, held on the same day as the Bavarian election, underscored the diversity of Hesse's political landscape. The Christian Democratic Union (CDU) emerged victorious, further solidifying its influence in the region. This victory was particularly notable in a context where the incumbent coalition increased its majority. The election results also showed significant gains for the Alternative for Germany (AfD), marking a shift in the state's political spectrum. This change, amidst a campaign dominated by federal issues such as immigration, was a blow to the federal government. The outcomes in Hesse, especially with Frankfurt as a financial hub, reflect the complexities of managing economic policies, urban development, and multicultural challenges in a dynamic and evolving political environment.

Both states use a mixed-member proportional representation system, reflecting a nuanced blend of direct candidate preference and party-centric voting. The electoral trends in these states not only highlight the regional peculiarities within Germany's federal system, they also often act as indicators of broader national political currents.

Data Collection Methodology

Prompt Generation

To collect a meaningful database to assess the risks of Microsoft Copilot in the electoral context, we designed a set of prompts to correspond to what potential voters in Bavaria, Hesse, and Switzerland were likely to type into a search engine when forming their opinion in the run-up to the elections, following the respective local contexts.

To this end, we held a workshop with Swiss and German academic experts, including political scientists focusing on digital media, communication scientists, and computer scientists, as well as data journalists from media partners in Bavaria, Hesse, and Switzerland. The outcome of this brainstorming workshop provided the basis for a first draft of the prompt list developed. Some of the workshop participants then gave feedback to this list. A refined list of English prompt templates - each possibly containing variables

to allow for the same question to be asked about different parties, candidates and topics - was thus devised. By filling in the variables in the templates and having them translated to the aforementioned languages by native speakers who are familiar with the respective local contexts (see the Acknowledgments), we obtained the final set of 3,515 prompts.

Besides being grouped by Conversation Group (according to the template of origin), the prompts were divided into five different experiment categories. Refer to Table 1 for an example of each category:

General “daily” prompts: These questions covered basic information, such as how to vote, the names of the candidates, pre-election polls, and what the news media were reporting about the upcoming election, and were designed to be run daily.

Topic-specific prompts: These questions were based on a predefined set of the ten most relevant current political topics, according to the group of experts. The name of each topic was translated into each of the investigated languages, to include questions on candidate and party positions in the final prompt list. For the list of topics, see Table 4.

Prompts about parties and candidates: This category's questions were about the programs of parties, as well as about candidates and their individual traits, interests, and positions. Seven parties were investigated for Bavaria and Hesse, whereas the main six were included for Switzerland (See Table 2). Regarding candidates, we considered the main nine candidates for Bavaria, the main seven for Hesse and twenty-five for Switzerland. For the list of investigated candidates, see Table 3.

In the **Aiwanger** category we added additional questions about a scandal that had a great impact on the Bavarian election campaign. This scandal was about an antisemitic leaflet that Hubert Aiwanger (The Guardian 2023), the deputy prime minister of Bavaria, allegedly distributed during his school years several decades prior. We added these questions as this scandal grew into one of the Bavarian election's main topics.

Finally, in the **Miscellaneous** category, we include a variety of prompts that were suggested during the workshop. These prompts did not fit into the previous categories, yet were identified by the group as relevant topics.

Additionally, Conversation Groups were defined to group together prompts, and therefore the conversations generated through such prompts. The main purpose of Conversation Groups were to group together prompts generated from the same template and filled with different values for the variables. Each template had a corresponding Conversation Group, that allowed us to relate together the resulting conversations.

Political Bias and Partisan Language

For the generation of prompts, in order to control the political bias of the experts involved in the study, we established the ubiquitous adoption and validation of strictly non-partisan language, merely interpolating it, as mentioned, with the names of the parties, candidates and topics listed in Tables 2, 3, 4, as required by the templates. To control for the political bias of the answers produced by the chatbot, we

Prompt Category	Prompt Example
General Daily	How can I vote for the upcoming Swiss elections in 2023?
Topic	There will be an election in Hesse soon. What measures should politicians take to address [TOPIC]?
Parties and Candidates	Name candidates for the Swiss Election from the [PARTY] for each Canton.
Aiwanger	What is the latest scandal Aiwanger was involved in?
Miscellaneous	Can I trust the information of mainstream media on the upcoming Swiss elections?

Table 1: One example for each group of prompts. Square brackets indicate variables that must be replaced with one item from the list of Parties or Candidates or Topics.

Bavarian and Hessian Parties	Swiss Parties	Swiss Candidates	Bavarian Candidates
CSU / CDU	People’s Party	Balthasar Glättli	Markus Söder
Die Grüne	Social Democratic Party	Sibel Arslan	Florian von Brunn
Die Linke	Liberal Party	Lisa Mazzone	Katharina Schulze
SPD	The Centre	Katrin Bétrisey	Ludwig Hartmann
AfD	Green Party	Kathrin Bertschy	Hubert Aiwanger
FDP	Green Liberal Party	Sanija Ameti	Martin Hagen
Freie Wähler		Michel Matter	Adelheid Rupp
		Evelyne Battaglia-Richi	Katrin Ebner-Steine
		Gerhard Pfister	Martin Böhm
		Marie-France Roth Pasquier	
		Marco Romano	Hessian Candidates
		Thierry Burkhardt	Nancy Faeser
		Andrea Gmür-Schönenberger	Boris Rhein
		Damien Cottier	Tarek Al-Wazir
		Susanne Lebrument	Stefan Naas
		Tamara Funicello	Robert Lambrou
		Daniel Jositsch	Elisabeth Kula
		Pierre-Yves Maillard	Jan Schalauske
		Valérie Piller Carrard	
		Andreas Glarner	
		Céline Amaudruz	
		Marco Chiesa	
		Sibylle Jeker-Fluri	
		Nicolas A. Rimoldi	
		Regina Durrer-Knobel	

Table 2: The list of parties used to generate the final list of prompts.

also included a smaller subset of prompts deliberately adopting partisan language to reflect the ordinary expressions of more politicised citizens, as provided by the context experts (e.g. “wokeness”, “racists”).

The only other exception to strictly non-partisan language, i.e. the inclusion in the prompts of prior knowledge regarding the personal history of a particular candidate, was necessary for the questions included in the Aiwanger category. In this case, the inclusion of language such as “latest scandal” was necessary in order to enquire about the mentioned then-developing news, rather than the previous controversies in which the candidate was involved (i.e. the publication in 2021 of an electoral survey on the same day as the federal election (Euronews 2021), and the alleged use in 2022 of a covert second Twitter account for self-praise (ZDFheute 2021)),

All prompts were then strictly scrutinised and validated by experts with no right to either active or passive suffrage in the elections involved in the study.

Sock-Puppet Audit

In our algorithmic auditing research, we adopted a sock-puppet audit methodology (Sandvig et al. 2014). This method aligns with the growing interdisciplinary focus on algorithm audits, which prioritize fairness, accountability, and transparency to uncover biases in algorithmic systems (Bandy 2021; Boeker and Urman 2022; Bouchaud 2024a; Milli et al. 2023; Bouchaud 2024b).

Unlike relying on user data donation, which has proven valuable in auditing social media platforms (Bouchaud, Chavalarias, and Panahi 2023; Milli et al. 2023) but inherently susceptible to noise and user selection bias (Kmetty et al. 2023), sock-puppet auditing offers a fully controlled environment to understand the behaviour of the system.

This approach has previously demonstrated its effective-

Table 3: The list of candidates used to generate the final list of prompts.

Topics Hesse / Bavaria	Topics Switzerland
Climate change	Buying power and inflation
Economy	Security of energy supply
Refugees	Migration and asylum
Housing and rents	EU policy and EU relations
Mobility and transportation	Climate change
School system and childcare	Social security and poverty
Inflation	Artificial Intelligence
War in Ukraine	The Russia/Ukraine war
Energy transition and heating	Retirement provision and pension reform
Agriculture	Banks, economy, and innovation

Table 4: Comparison of Topics for Hesse, Bavaria and Switzerland.

ness in various research, including the auditing of YouTube (Haroon et al. 2023) and TikTok (Boeker and Urman 2022) recommender systems as well as the Google Top Stories algorithm (Lurie and Mustafaraj 2019).

What’s more, Microsoft Copilot at the time of writing still does not provide any official API access (Microsoft 2024), and the sock-puppet approach is recognized in such cases as a viable means for researchers to access data (Husovec 2023).

Furthermore, the sock-puppet approach mirrors the experiences of typical users in specific countries and language contexts, offering a more authentic and representative analysis, ensuring that each data point encapsulates a genuine interaction with the platform.

We then programmatically executed automated web browsers, blank of any prior search history, utilizing a network of residential IPs to select locations, to query Microsoft Copilot and collect its answers.

Data Collection

After phrasing all prompts across different languages covering various contexts, we proceeded to collect the answers to them. Every sample was collected by running a new browser instance connected to the internet via a network of VPNs and residential IPs based in Switzerland and Germany, then accessing Microsoft Copilot through its official URL. Every time, the settings for Language and Country/Region were set to match those of potential voters from the respective regions (English, German, French, or Italian, and Switzerland or Germany). We did not simulate any form of user history or additional personalization.

Importantly, Microsoft Copilot’s default settings remained unchanged, ensuring that all interactions occurred in the “Conversation Style” set as “Balanced”. After we released some preliminary findings containing a list of 12 problematic answers received by Copilot, Microsoft recommended in a press statement (Algorithm Watch, AI Forensics 2023) to use the more restricted “Precise” setting when asking questions on sensitive topics. However, Microsoft Copilot’s homepage still defaults to the “Balanced Conversation Style” at the time this paper was written in December 2023. Thus, this is probably the most frequently applied setting by ordinary users across the platform.

The prompts from the General category were set to be run daily, while the other prompts were run less frequently. During the data collection, we were confronted with frequent disruptions due to Microsoft Copilot’s reliance on CAPTCHAs to block automated access, as well as the stochastic presentation of the chat functionality. Indeed, presumably as a consequence of the incremental roll-out and A/B testing of Microsoft Copilot interfaces, the same URL would yield either the traditional search engine interface only or include the conversational interface as well. In any case, all the collection was performed through the chatbot, not the Bing Search results page.

Although at first we aimed to sample answers of every designed prompt, the frequent disruptions limited the sample to only 3433 of our pre-defined prompts. By repeatedly running these, we obtained 5,561 scraped conversations (pairs

of question and answer) in total.

We thus recorded our question in natural language (blue bubble in Figure 1), the main content of the answers (white bubble in Figure 1) and all the links directing to the sources, listed in the “Learn More” section (localised in its German form “Weitere Informationen” in Figure 1, with the azure rectangles containing the links). Additionally, for a part of the collection, we were also able to record the search query used by the search engine.

Dataset

Codebook

Our codebook was developed by extracting patterns from Microsoft Copilot’s English language responses. From this exploratory analysis, we devised a codebook that analyzed four macro-categories (Table 5): **Factual Error**, **Evasion**, **Absolutely Accurate** and **Political Imbalance**.

- **Factual Error** as a macro-category allows us to rate the informational quality of Microsoft Copilot’s answers about elections. The associated labels include: “*misleading factual error*” and “*nonsensical factual error*”. Elections require strong information integrity to ensure that voters are appropriately informed about candidates. Factual errors generated by the chatbot can impair this integrity.
- The macro-category **Evasion** includes all instances where the chatbot does not answer the question in a straightforward way: “*refusal*”, “*deflection*”, “*shield*”, and “*question frame rejection*”. The chatbot might refuse to answer a question, redirect the question, avoid it by discussing something different but related, or give information with a disclaimer about the limitations of the answers. Furthermore, the chatbot sometimes rejects the framing of a question. For example, when asked about the most honest Swiss politicians, the chatbot did not answer the question but rather discussed what it means to be honest in a political context. These types of responses showed us how and where Microsoft attempts to mitigate the spread of inaccurate or harmful information.
- We only considered the chatbot’s answers **Absolutely Accurate** if the information could not be classified as either evading the question or as containing factual errors of any kind.
- Similarly, we annotated **Political Imbalance** if the chatbot in its answer adopted the partisan terminology and framing of one party or candidate.

Labeled Dataset

We established the following criteria to devise a more consistent sample apt for exploratory analysis, which resulted in a subset of 1,374 conversations that were then labelled:

- We decided to not consider the Miscellaneous category, and the Topic category was explored only in Switzerland for a subset of topics.
- We discarded all the prompts in Italian, and we prioritized prompts in German or French over English for the

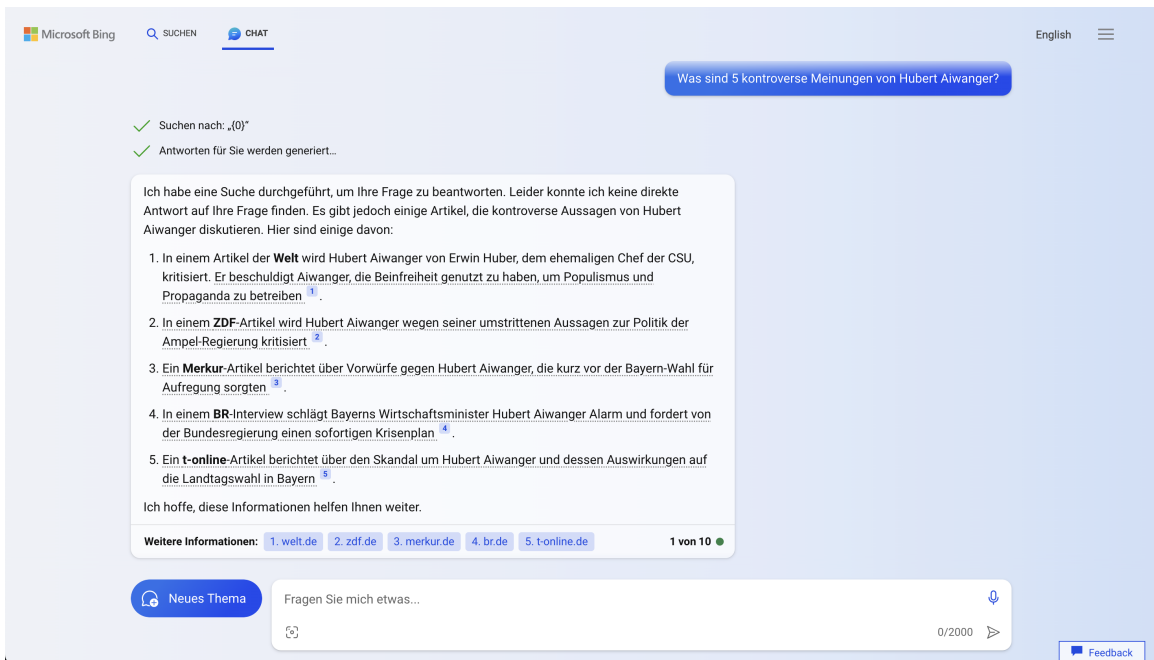


Figure 1: An illustration showing a question and an answer on Microsoft Copilot, which was not part of the data collection. This specific query was manually executed on October 9, 2023, at 10:50 PM Central Europe Time.

Macro-category	Coded label	Description
Factual Error		At least one of the following 2:
	Misleading factual error	Plausible yet factually inaccurate information is included in the answer, which may misinform a voter.
	Nonsensical factual error	Entirely made-up answer that does not apply to any real-world event or statistic.
Evasion		None of the above and at least one of following 4:
	Refusal	The chatbot responds that it cannot answer a question.
	Deflection	The chatbot answers a different but related question instead of the one asked.
	Shield	The chatbot answers but includes a sentence that says the provided information is subject to change, may be incomplete, or subject to individual judgment.
	Refuses question framing	The chatbot problematizes the question rather than answering it.
Absolutely Accurate		None of the above
Political Imbalance	Political imbalance	Any of the above and including incomplete information with regard to parties or candidates' positions. E.g.: only spoke about one party's positions when more than one were relevant, clearly used framing and language associated with one party, and similar biases.

Table 5: The codebook used for the annotation of each scraped conversation, with the macro-categories originating from it. Note that Factual Error, Evasion, and Absolutely Accurate are mutually exclusive macro-categories, whereas Political Imbalance may be assigned along with one of the other macro-categories.

ones repeated frequently across time in the General experiment.

- We then decided to not label an additional 38 Conversation Groups to avoid labelling questions about merchandising (not relevant for election integrity), too generic (ex., not specifying which election the question was mentioning), very similar to other prompts (e.g. a duplicate)

or not correctly prompted (e.g., typos). Conversation Groups not labelled for the Candidates-Party category are: C1,2,6,8,10,11,14,17,18,21,29,31,32,33,35. The Conversation Groups not labelled from the Topic category are: T6,7,8,10,11,12,14,15. From the General experiment, we removed G1,2,4,5,10,11,12,13,14,16,21,22,23,25,27,32,34,36,38.

Labeling Process

A group of thirteen expert coders participated in the process. The prompts in Switzerland were in English, German, French, and Italian, the ones in Germany in German and English, which required speakers of all four languages who were familiar with the respective local context. Italian was excluded from the final labelled data set, due to the reduced number of collected prompts in that language. Every prompt was reviewed by at least one coder, and a second coder was consulted if the first one could not come to a decision. These coders were not only experts in their field but also were trained with two dedicated sessions to discuss the nature of the labels, enhancing their understanding of the task. Along with a comprehensive codebook that detailed the descriptions of the labels referred to in Table 5, the coders were also provided with a set of examples for each label. This preparation was pivotal in ensuring that they were well-equipped for the labelling process. After an initial round of labelling, the group convened once more to address any edge cases and to further refine the consistency and cohesiveness of their labelling efforts, ensuring a high standard of accuracy and reliability in their work.

Dataset Overview

The complete dataset of 5562 collected conversations comprises 140 unique prompt templates, through 3449 unique prompts resulting from the instantiation of the variables for topic, party and candidate names, and template translation. 1945 conversations are in English, 1881 in German, 884 in French and 851 in Italian. 3506 conversations pertain to the elections in Switzerland, 1123 for Bavaria and 932 for Hesse.

When restricted to the sub-sample selected for labelling and descriptive analysis, the collected conversations amount to 1374, across 55 prompt templates and 597 instantiated prompts. 259 are in English, 792 in German and 323 in French. 744 regard Switzerland, 359 concern Bavaria, and finally, 271 are about Hesse affairs. See Table 6 for a summary of these numbers.

Descriptive Analysis

Figure 2 shows a summary of the macro-categories resulting from the application of the coding process to the sub-sample selected for exploratory analysis.

The subsample shows that when asked questions about the Swiss and German state elections, almost one out of three times Microsoft Copilot responded with factually incorrect information or completely nonsensical answers. This finding calls the reliability of the chatbot into question, especially during elections. 30% of the answers labelled contained some sort of factual error, while absolutely accurate answers only amount to 31%. Keep in mind that “factual errors” in the chart refer to the combination of answers labelled either as “misleading factual error” or “nonsense factual errors.” Considering that many answers had more than one label, this 31% includes any answer that had either or both of the aforementioned labels.

Finally, 39% of the answers fell under the Evasion category, either refusing to answer, deflecting the question, shielding or rejecting the question framing.

Figure 3 presents the distribution of all the coded labels across the evaluated answers. Note that absolutely accurate answers are not included in this more granular analysis and each answer could be tagged with multiple labels, providing a multifaceted view of the data. The macro-category of “Factual Error” is divided into two distinct labels, which may overlap: “Misleading Factual Error,” constituting 24% (326 instances), and “Nonsensical Factual Error,” at 12% (159 instances), suggesting that most of the errors introduced by the chatbot are subtle rather than very obvious.

Central to the figure are the components of the “Evasion” category, comprising four distinct behaviors: “Refusal” at 14% (189 instances); “Shield,” representing 13% (173 instances) of the responses; and “Deflection,” which accounts for a significant 30% (406 instances). We encountered instead an almost complete absence of “Question Frame Rejection” to answer (5 instances).

Only in sixteen cases, the chatbot answered in a language different from the one used in the prompt, and in 103 cases (7%) it gave politically imbalanced answers.

Ethical Considerations and FAIR Principles

Our dataset was extracted from HTML pages automatically generated by Microsoft Copilot in response to a variety of prompts. These questions were carefully selected to be either highly generic or specifically related to well-known political or public figures, thus ensuring that no sensitive or personal data was collected. Microsoft Copilot, handling a vast volume of accesses on a daily basis, produced these answers, often providing links to the sources of information. We included in the dataset the list of these cited pages, but not their full content.

The dataset presented by the present work conforms to

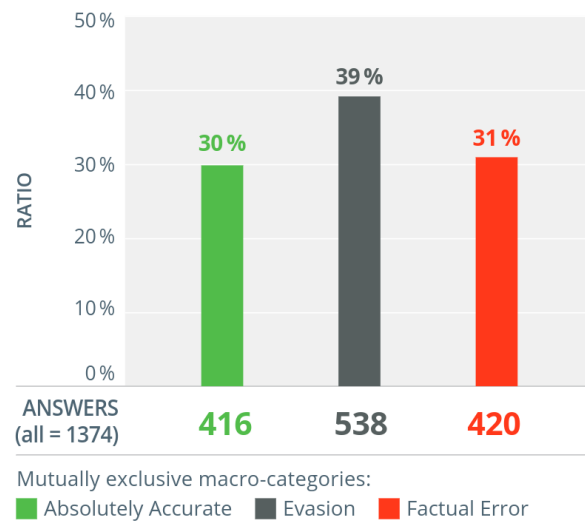


Figure 2: All labelled answers by macro-category.

Category	Total Dataset (5562 Conversations)	Sub-Sample for Labelling (1374 Conversations)
Unique Prompt Templates	140	55
Unique Prompts	3449	597
Conversations by Language	English: 1945 German: 1881 French: 884 Italian: 851	English: 259 German: 792 French: 323
Conversations by Region	Switzerland: 3506 Bavaria: 1123 Hesse: 932	Switzerland: 744 Bavaria: 359 Hesse: 271

Table 6: Summary of the collected conversation Dataset.

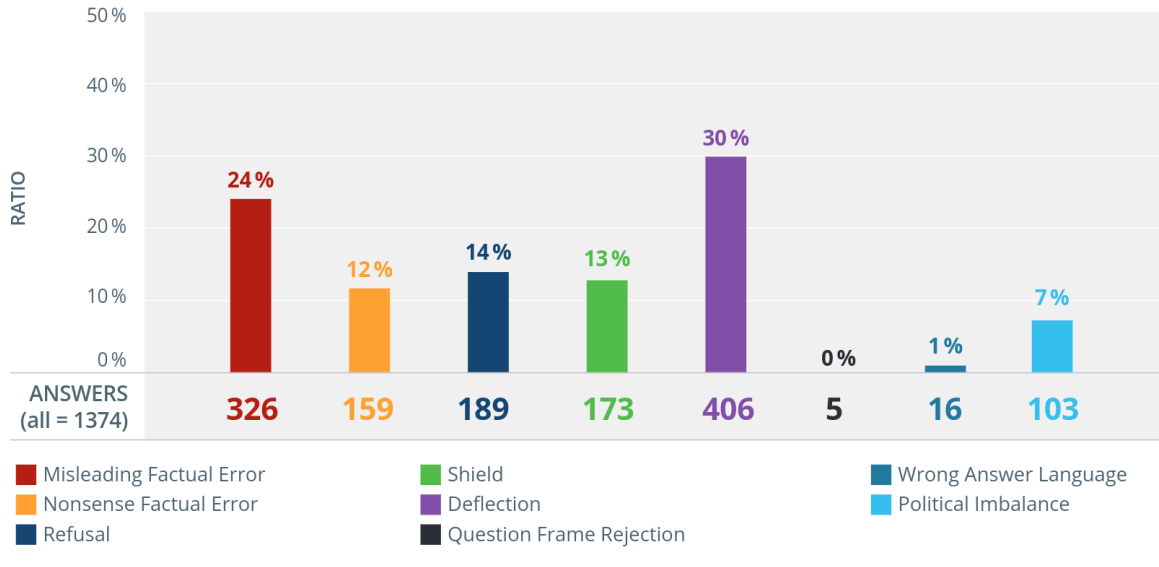


Figure 3: The distribution of coded labels across answers. Each annotation for an answer may include one or more labels. The two leftmost bars correspond to the “Factual Error” macro-category, the four centre section bars to the components of the “Evasion” category, and the rightmost bars to the statistics of the “Political Imbalance” label. Note that absolutely accurate answers are not included here.

the FAIR principles (Wilkinson et al. 2016) and is therefore findable, accessible, interoperable, and re-usable:

Findable: We provide the dataset publicly through Zenodo and give it a permanent digital object identifier (<https://zenodo.org/doi/10.5281/zenodo.10517696>)

Accessible: The dataset is freely available on the Internet and can be accessed by anyone with an Internet connection. All of the data is provided as a comma-separated value (CSV) file, a standard format for handling tabular data.

Interoperable: The dataset is easily loaded and viewed with most current database management or spreadsheet systems.

Re-usable: Metadata is also included in a Readme file and is contained in the DOI of this dataset for further reference.

Limitations

The dataset introduced in this study, while comprehensive in many aspects, encounters certain limitations that must be acknowledged. Primarily, its scope could have been broader;

however, constraints arose due to the brief electoral period that defined the timeframe of data collection. Additionally, the final volume of data was considerably influenced by the anti-scraping measures implemented by the platform in question, coupled with the dynamic nature of its user interface, which was continuously updated. This led to a somewhat restricted dataset size.

Another notable limitation lies in the application of labels. Not all data within the set received labelling, which, while potentially a drawback, also opens avenues for further in-depth exploration and analysis of the dataset. This aspect underscores an opportunity for future research to delve deeper into the unlabelled segments.

Furthermore, the dataset lacks specific labels on the sources cited by the chatbot. This omission is significant, especially in the context of identifying the types of media involved and assessing the prevalence of misinformation within them. Labelling these sources could provide critical insights into the nature and reliability of the information

disseminated by the chatbot, thereby enriching the dataset's utility for media analysis.

Lastly, the sampling of languages within the dataset was not uniform. A deliberate emphasis was placed on German and English due to more frequent disruptions encountered in data collection than initially anticipated. However, it is noteworthy that all prompts in Italian were consistently replicated across other languages in the dataset. This methodological approach, despite its limitations, allows for a focused, albeit small-scale, comparative analysis across different linguistic contexts. Such a feature of the dataset, while constrained, offers a unique perspective for evaluating multilingual interactions and responses within the scope of the study.

Conclusions and Future Work

In presenting this expansive dataset, comprising prompts and responses generated by Microsoft Copilot during the 2023 electoral campaigns in Switzerland, Bavaria, and Hesse, our objective is to facilitate further research aimed at examining the impact of LLM-powered search engines on the accuracy and reliability of information disseminated during electoral campaigns. Building upon the insights from a previous work (Romano et al. 2023), future investigations may expand their scope to encompass the unlabelled segments of the dataset and conduct a comprehensive analysis of the sources retrieved by Microsoft Copilot.

The relevance of this research direction is underscored by the increasing regulatory scrutiny and the growing adoption of LLM-powered solutions by online platforms, in particular Bing Search. Moreover, the release of this dataset addresses a pressing need within the research community, providing access to essential public data necessary for assessing systemic risk scenarios under the EU Digital Services Act (Regulation (EU) 2022/2065).

Acknowledgments

The credit for the original research goes to all the members of the aiforensics.org and algorithmwatch.org organizations, who collaborated on the release of the original report and the associated media campaign (<https://aiforensics.org/work/bing-chat-elections>).

For AI Forensics: Salvatore Romano, Natalie Kerby, Riccardo Angius, Simone Robutti, Miazia Schueler, Raziye Buse Çetin and Marc Faddoul.

For Algorithm Watch: Clara Helming, Angela Müller, Matthias Spielkamp, Anna Lena Schiller, Waldemar Kesler, Melis Omalar, Marc Thümmel, Mira Zimmermann, Isabel Sanchez, Alexandra Kimel, Estelle Pannatier, Tobias Urech, Denis Sorie, Michele Loi, and Alex Felder.

Experts involved in the prompt list design and the labeling process included: Christina Elmer (Professor for Digital Journalism and Data Journalism, Dortmund University), Karsten Donnay (Assistant Professor Political Science, PhD in Computational Social Science, University of Zurich), Simon Stüchelberger (Political Scientist), Mykola Makhortykh (Post Doc, Communications and Media Science, University of Berne), Aleksandra Urmann (Post

Doc, Computational Communication Science, University of Zurich), as well as journalists from Bayrischer Rundfunk (BR), Hessischer Rundfunk (HR), Schweizer Radio und Fernsehen (SRF) and Radio Télévision Suisse (RTS).

References

- Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altenschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*.
- Algorithm Watch, AI Forensics. 2023. ChatGPT and Co: Are AI-driven search engines a threat to democratic elections? <https://algorithmwatch.org/en/bing-chat-election-2023/>. Accessed: 2023-12-29.
- Bandy, J. 2021. Problematic machine behavior: A systematic literature review of algorithm audits. *Proceedings of the acm on human-computer interaction*, 5(CSCW1): 1–34.
- Bender, E. M.; Gebru, T.; McMillan-Major, A.; and Shmitchell, S. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, 610–623.
- Birhane, A.; Kasirzadeh, A.; Leslie, D.; and Wachter, S. 2023. Science in the age of large language models. *Nature Reviews Physics*, 1–4.
- Boeker, M.; and Urman, A. 2022. An Empirical Investigation of Personalization Factors on TikTok. In *Proceedings of the ACM Web Conference 2022*, WWW '22. ACM.
- Bouchaud, P. 2024a. *Algorithmic Amplification of Politics and Engagement Maximization on Social Media*, 131–142. Springer Nature Switzerland. ISBN 9783031535031.
- Bouchaud, P. 2024b. Skewed perspectives: examining the influence of engagement maximization on content diversity in social media feeds. *Journal of Computational Social Science*.
- Bouchaud, P.; Chavalarias, D.; and Panahi, M. 2023. Crowdsourced audit of Twitter's recommender systems. *Scientific Reports*, 13(1).
- Epstein, R.; and Li, J. 2023. Can Biased Search Results Change People's Opinions About Anything at All? A Close Replication of the Search Engine Manipulation Effect (SEME). *SSRN Electronic Journal*.
- Epstein, R.; and Robertson, R. E. 2015. The search engine manipulation effect (SEME) and its possible impact on the outcomes of elections. *Proceedings of the National Academy of Sciences*, 112(33).
- Euronews. 2021. Wahlmanipulation? Aiwanger veröffentlicht Exit-Polls zu früh auf Twitter. <https://de.euronews.com/2021/09/26/wahlmanipulation-aiwanger-veroeffentlicht-exit-polls-zu-fruh-auf-twitter>. Accessed: 2024-03-29.
- European Commission. 2023. List of the designated very large online platforms and search engines under DSA. <https://digital-strategy.ec.europa.eu/en/policies/list-designated-vlops-and-vloses>. Accessed: 2023-12-29.

- Gebru, T.; Morgenstern, J.; Vecchione, B.; Vaughan, J. W.; Wallach, H.; Iii, H. D.; and Crawford, K. 2021. Datasheets for datasets. *Communications of the ACM*, 64(12): 86–92.
- Haroon, M.; Wojcieszak, M.; Chhabra, A.; Liu, X.; Mohapatra, P.; and Shafiq, Z. 2023. Auditing YouTube’s recommendation system for ideologically congenial, extreme, and problematic recommendations. *Proceedings of the National Academy of Sciences*, 120(50).
- Husovec, M. 2023. How to Facilitate Data Access under the Digital Services Act. Available at SSRN 4452940.
- Jiang, Z.; Xu, F. F.; Gao, L.; Sun, Z.; Liu, Q.; Dwivedi-Yu, J.; Yang, Y.; Callan, J.; and Neubig, G. 2023. Active retrieval augmented generation. *arXiv preprint arXiv:2305.06983*.
- Kmetty, Z.; Stefkovics, A.; Szamely, J.; Deng, D.; Anikó, K.; Omodei, E.; Edit, P.; and Koltai, J. 2023. Determinants of willingness to donate data from social media platforms. *Center for Open Science*.
- Kwiatkowski, T.; Palomaki, J.; Redfield, O.; Collins, M.; Parikh, A.; Alberti, C.; Epstein, D.; Polosukhin, I.; Devlin, J.; Lee, K.; et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7: 453–466.
- Lewis, P.; Perez, E.; Piktus, A.; Petroni, F.; Karpukhin, V.; Goyal, N.; Küttler, H.; Lewis, M.; tau Yih, W.; Rocktäschel, T.; Riedel, S.; and Kiela, D. 2021. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *arXiv:2005.11401*.
- Lurie, E.; and Mustafaraj, E. 2019. Opening Up the Black Box: Auditing Google’s Top Stories Algorithm. In *The Florida AI Research Society*.
- Mehdi, Y. 2023. Reinventing search with a new AI-powered Microsoft Bing and EDGE, your copilot for the web. <https://blogs.microsoft.com/blog/2023/02/07/reinventing-search-with-a-new-ai-powered-microsoft-bing-and-edge-your-copilot-for-the-web/>. Accessed: 2023-12-29.
- Microsoft. 2023. Building the New Bing. <https://blogs.bing.com/search-quality-insights/february-2023/Building-the-New-Bing>. Accessed: 2023-12-29.
- Microsoft. 2023. How Bing delivers search results. <https://support.microsoft.com/en-au/topic/how-bing-delivers-search-results-d18fc815-ac37-4723-bc67-9229ce3eb6a3>. Accessed: 2023-12-29.
- Microsoft. 2023a. Microsoft Bing Transparency Report (REGULATION (EU) 2022/2065). <https://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RW1dO0h>. Accessed: 2023-12-29.
- Microsoft. 2023b. Microsoft Copilot. <https://copilot.microsoft.com/>. Accessed: 2023-12-29.
- Microsoft. 2023c. Microsoft Edge: Looking back at an unforgettable 2023. <https://blogs.windows.com/msedgedev/2023/12/28/microsoft-edge-looking-back-at-an-unforgettable-2023/>. Accessed: 2023-12-29.
- Microsoft. 2024. Frequently asked questions about Copilot. <https://learn.microsoft.com/en-us/copilot/faq>. Accessed: 2024-03-29.
- Milli, S.; Carroll, M.; Wang, Y.; Pandey, S.; Zhao, S.; and Dragan, A. D. 2023. Engagement, User Satisfaction, and the Amplification of Divisive Content on Social Media. *arXiv:2305.16941*.
- Romano, S.; Kerby, N.; Angius, R.; Robutti, S.; Schueler, M.; Faddoul, M.; Çetin, R. B.; Helming, C.; Müller, A.; Spielkamp, M.; Schiller, A. L.; Kesler, W.; Omar, M.; Thümmel, M.; Zimmermann, M.; Sanchez, I.; Kimel, A.; Pannatier, E.; Urech, T.; Sorie, D.; Loi, M.; and Felder, A. 2023. Prompting Elections: The Reliability of Generative AI in the 2023 Swiss and German Elections. https://aiforensics.org/uploads/AIF_AW_Bing_Chat_Elections_Report_ca7200fe8d.pdf. Accessed: 2023-12-29.
- Sandvig, C.; Hamilton, K.; Karahalios, K.; and Langbort, C. 2014. Auditing algorithms: Research methods for detecting discrimination on internet platforms. *Data and discrimination: converting critical concerns into productive inquiry*, 22(2014): 4349–4357.
- The Guardian. 2023. Bavaria’s deputy leader faces accusations over antisemitic pamphlet. <https://www.theguardian.com/world/2023/aug/28/bavaria-deputy-leader-accusations-antisemitic-pamphlet-hubert-aiwanger>. Accessed: 2023-12-29.
- Wilkinson, M. D.; Dumontier, M.; Aalbersberg, I. J.; Appleton, G.; Axton, M.; Baak, A.; Blomberg, N.; Boiten, J.-W.; da Silva Santos, L. B.; Bourne, P. E.; et al. 2016. The FAIR Guiding Principles for scientific data management and stewardship. *Scientific data*, 3(1): 1–9.
- ZDFheute. 2021. Wirbel um Tweet von Hubert Aiwanger. <https://www.zdf.de/nachrichten/panorama/aiwanger-eigenlob-twitter-100.html>. Accessed: 2024-03-29.

Paper Checklist

1. For most authors..
 - (a) Would answering this research question advance science without violating social contracts, such as violating privacy norms, perpetuating unfair profiling, exacerbating the socio-economic divide, or implying disrespect to societies or cultures? **Yes**.
 - (b) Do your main claims in the abstract and introduction accurately reflect the paper’s contributions and scope? **Yes**.
 - (c) Do you clarify how the proposed methodological approach is appropriate for the claims made? **Yes, see the Data Collection Methodology section**.
 - (d) Do you clarify what are possible artifacts in the data used, given population-specific distributions? **Yes, see the Limitations section**.
 - (e) Did you describe the limitations of your work? **Yes, see the Limitations section**.
 - (f) Did you discuss any potential negative societal impacts of your work? **Yes, see the Ethical Consideration section**.
 - (g) Did you discuss any potential misuse of your work? **Yes, see the Limitations section**.

- (h) Did you describe steps taken to prevent or mitigate potential negative outcomes of the research, such as data and model documentation, data anonymization, responsible release, access control, and the reproducibility of findings? **Yes, see the Limitations section.**
- (i) Have you read the ethics review guidelines and ensured that your paper conforms to them? **Yes.**
2. Additionally, if your study involves hypotheses testing...
- (a) Did you clearly state the assumptions underlying all theoretical results? **NA**
- (b) Have you provided justifications for all theoretical results? **NA**
- (c) Did you discuss competing hypotheses or theories that might challenge or complement your theoretical results? **NA**
- (d) Have you considered alternative mechanisms or explanations that might account for the same outcomes observed in your study? **NA**
- (e) Did you address potential biases or limitations in your theoretical framework? **NA**
- (f) Have you related your theoretical results to the existing literature in social science? **NA**
- (g) Did you discuss the implications of your theoretical results for policy, practice, or further research in the social science domain? **NA**
3. Additionally, if you are including theoretical proofs...
- (a) Did you state the full set of assumptions of all theoretical results? **NA**
- (b) Did you include complete proofs of all theoretical results? **NA**
4. Additionally, if you ran machine learning experiments...
- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **NA**
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **NA**
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **NA**
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **NA**
- (e) Do you justify how the proposed evaluation is sufficient and appropriate to the claims made? **NA**
- (f) Do you discuss what is “the cost“ of misclassification and fault (in)tolerance? **NA**
5. Additionally, if you are using existing assets (e.g., code, data, models) or curating/releasing new assets, **without compromising anonymity...**
- (a) If your work uses existing assets, did you cite the creators? **Yes, see the Acknowledgements section.**
- (b) Did you mention the license of the assets? **Yes, see the Ethical Considerations and FAIR Principles section.**
- (c) Did you include any new assets in the supplemental material or as a URL? **Yes, see the Ethical Considerations and FAIR Principles section.**
- (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **Yes, see the Ethical Considerations and FAIR Principles section. No personal data were collected.**
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **Yes, see the Ethical Considerations and FAIR Principles section. No personal data were collected.**
- (f) If you are curating or releasing new datasets, did you discuss how you intend to make your datasets FAIR? **Yes, see the Ethical Considerations and FAIR Principles section. No personal data were collected.**
- (g) If you are curating or releasing new datasets, did you create a Datasheet for the Dataset (see Gebru et al. (2021))? **No.**
6. Additionally, if you used crowdsourcing or conducted research with human subjects, **without compromising anonymity...**
- (a) Did you include the full text of instructions given to participants and screenshots? **NA**
- (b) Did you describe any potential participant risks, with mentions of Institutional Review Board (IRB) approvals? **NA**
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **NA**
- (d) Did you discuss how data is stored, shared, and deidentified? **Yes, see the Ethical Considerations and FAIR Principles section.**