

---

# The Tournesol dataset: Which videos should be more largely recommended?

---

Anonymous Author(s)  
Affiliation  
Address  
email

## Abstract

1 This paper introduces the Tournesol public dataset, which was collected as part  
2 of the online deployed platform <https://tournesol.app>. Our dataset contains  
3 a list of 204,000 comparative judgments made by Tournesol’s 20,000 users on  
4 which YouTube videos should be more largely recommended. It also provides  
5 703,000 comparisons along secondary criteria like content reliability, topic impor-  
6 tance and layman-friendliness. The dataset also exports information about users’  
7 pretrust statuses and vouches. It is published at [https://api.tournesol.app/](https://api.tournesol.app/exports/all)  
8 [exports/all](https://api.tournesol.app/exports/all) under ODC-By license. The data is currently used by Tournesol to  
9 make community-driven video content recommendations to over 10,000 users.

## 10 1 Introduction

11 Recommendation algorithms have become extremely influential. In the last few years, beyond their  
12 impacts on mental health [54, 19, 91], because they amplify disinformation, cyberbullying and hate,  
13 they have been linked to major geopolitical events, including COVID disinformation [78, 43], the rise  
14 of far-right parties [90, 89, 94], and the Rohingya genocides [39, 71]. Crucially, in all these examples,  
15 the victims of recommendation algorithms are not only their users; hate amplification is threatening  
16 entire populations, even when these populations do not use recommendation algorithms themselves.  
17 This is in sharp contrast with the overwhelming majority of the scientific literature, which assumes  
18 that recommendation algorithms should be optimized for their users only [1, 69].

19 As online activities grew, social media have *de facto* taken the role that was traditionally played by  
20 these intermediate bodies [88, 47]. This became particularly striking when, in 2020, the then US  
21 President was banned from Twitter, Facebook, and Youtube, long before any court sued him for  
22 inciting the Capitol riot violence [64, 65]. As another example, by amplifying the cyberbullying of  
23 climate scientists, Twitter provoked their exodus from the platform [92], thereby turning climate  
24 change into a *mute news*, which is endangering plenty of non-users [3]. The great replacement of the  
25 intermediate body by privately owned algorithms has been tied to an alarming decline of democratic  
26 norms worldwide, as many reports expose a global trend of autocratization [70, 7].

27 So how do today’s large-scale recommendation algorithms address the ethical dilemmas that they face  
28 billions of times per day, when they are tasked with amplifying some (potentially hateful) content over  
29 others (of potential public interest)? Currently, they heavily rely on (highly sophisticated) *machine*  
30 *learning* [23, 61]. In other words, such algorithms leverage massive amounts of data to determine  
31 which content they will promote at scale. However, as an immediate corollary, such algorithms are  
32 exposed to *manipulation* by *poisoning* data [86]. In fact, this poisoning has been industrialized, not  
33 only by authoritarian states [18, 45], but also by private companies based in the UK [49], Spain [14],  
34 Israel [6], France [87] and Switzerland [34]. The magnitude of this industry is well captured by one  
35 puzzling statistic: Facebook reportedly removes around *7 billion fake accounts per year* [56].

36 While a recent line of research has provided numerous poisoning mitigations [13, 31, 32, 27, 80, 74],  
37 it is also known that there are fundamental impossibility theorems that prevent accurate learning in  
38 highly adversarial, heterogeneous and high-dimensional settings [28, 57, 36, 30]. In particular, there  
39 is no substitute for training datasets of high quality and security. In particular, to design trustworthy  
40 ethical algorithms, it is essential to train them on large, secured and trustworthy datasets of human  
41 ethical judgments. In this paper, we present the *Tournesol public dataset*, whose goal is to remedy  
42 the current state of affairs. More precisely we make the following contributions.

43 **Contributions.** Our main contribution is to present and share the *Tournesol public dataset*, which  
44 can be downloaded directly from <https://api.tournesol.app/exports/all>. The dataset con-  
45 sists of over 204,000 pairwise comparisons of the recommendability of over 40,000 YouTube video by  
46 over 20,000 Tournesol accounts. Additionally, the dataset contains over 703,000 pairwise comparisons  
47 of the videos’ quality on secondary criteria, such as reliability, importance and layman-friendliness.  
48 Our dataset, published under ODC-By license, also contains pretrust information about contributors,  
49 vouches between contributors, as well as scores computed from the data using SOLIDAGO [12].  
50 Crucially, the dataset was collected in a fully deployed environment with actual stakes, as Tournesol  
51 eventually makes recommendations based on the provided data to over 10,000 users.

52 The paper also presents an analysis of our dataset, with valuable insights for the ethics of content  
53 recommendation. One finding is that the topic importance highly matters in Tournesol’s contributors’  
54 judgments. While caveats apply, this suggests that the attention to “fake news” may be misguided;  
55 in fact, the disinformation industry often proceeds *without* producing false information, e.g. by  
56 overclaiming positive impacts, shifting blame or bullying critics [75]. Prioritizing greater exposure  
57 to *mute news* might be more urgent. Our analysis also highlights the need of psychological-based  
58 preference learning models, as we expose biases and variations in contributors’ judgments.

59 Finally, our paper discusses numerous exciting research directions that our public dataset could  
60 inspire or facilitate. In particular, we believe that a lot more focus should be given to secure learning  
61 under poisoning attacks, but also to *Proof of Personhood*, *expertise validation*, *volition learning*,  
62 *active learning* and *resilient collaborative filtering*, among others.

63 **Literature review.** Tournesol presents a new contribution to the growing field of AI alignment with  
64 human values [46, 21, 50, 76], which aims to teach human preferences to algorithms, and to design  
65 systems that maximize what humans prefer to maximize [81, 52]. Clearly, this requires finding out  
66 about humans’ judgments on how algorithms ought to behave. Unfortunately, so far, to the best of  
67 our knowledge and especially for the important case of recommendation algorithms, there have not  
68 been many secure, public and free-license datasets with such AI-safety-critical data.

69 To collect such data in a realistic setting, Tournesol’s dataset draws inspiration from several previous  
70 AI ethics solutions, which leveraged *collaborative governance* to address cases of conflictual human  
71 judgments. In particular, [60] introduced WeBuildAI, a framework where stakeholders of a food  
72 donation system could weigh in on the identity of the recipient of a donation. One challenge is that  
73 such decisions must be made every day; but stakeholders are not available every time a decision needs  
74 to be made. To account for their preferences, WeBuildAI asks stakeholders to either write down  
75 an algorithm that describes their preferences, or to provide judgments on generated food donation  
76 dilemmas. In the latter case, a learning model is then used to infer how the stakeholders would likely  
77 assess other dilemmas. In any case, an *algorithmic representative* is thereby constructed for each  
78 stakeholder; and the resulting decision will follow from a vote of the algorithmic representatives.  
79 Similar approaches were proposed for kidney donation [42] and for the “trolley dilemmas” [40] that  
80 autonomous cars could one day face [10, 73].

81 Perhaps most similar to our approach are Twitter’s *Community Notes* [95, 77], whose governance  
82 is intended to be fully community-driven. More specifically, the system allows a community of  
83 contributors to add a note to misleading tweets, e.g. to correct misinformation or to add context  
84 to prevent confusion. The contributors cannot only propose the note; they are also asked to assess  
85 other contributors’ notes. Notes that are judged helpful by a sufficiently large and diverse set of  
86 contributors are then published by the platform. The system is very transparent, and provides a lot of  
87 freely accessible data on human judgments<sup>1</sup>.

---

<sup>1</sup>The data can be downloaded here: <https://communitynotes.twitter.com/guide/en/under-the-hood/download-data>

88 **Structure of the paper.** In the sequel, Section 2 will present our public dataset, and the context in  
89 which the data was provided. Section 3 presents an analysis of our dataset. Section 4 then provides a  
90 list of research challenges that are raised by the dataset. Finally, Section 5 concludes.

## 91 2 The dataset

92 In this section, we describe our main contribution, namely the release of a new, scalable, secured and  
93 trustworthy database of reliable human judgments.

### 94 2.1 Raw data

95 **Pretrust.** To guarantee the security of our data, Tournesol aims to verify that every account is  
96 owned and controlled by a human, and that this human only owns and controls this single account  
97 on the platform. In other words, Tournesol aims to obtain a *Proof of Personhood* [15] to verify each  
98 active Tournesol account, and to thereby prevent *Sybil attacks* [25]. Unfortunately, there is currently  
99 no reliable and scalable solution for *Proof of Personhood*.

100 Today’s main solution is *email certification*. More precisely, when they create a Tournesol account,  
101 contributors are asked to validate, if possible, an email address from a trusted email domain. The list  
102 of trusted email domains is currently managed manually. An email domain will be considered trusted  
103 if it seems sufficiently unlikely that a large number of fake accounts can be created from this domain.

104 This excludes domains like @gmail.com and personal domains like @my-personal-website.com.  
105 The concern is not only that the domain will maliciously create a large number of fake accounts; it  
106 is also that they may be hacked by a malicious entity that will create such fake accounts. The list  
107 of trusted email domains is available at [https://tournesol.app/email\\_domains](https://tournesol.app/email_domains). It includes  
108 domains like @epfl.ch, @who.int and @rsf.org. 703 contributors are thereby authenticated.

109 Evidently, however, this solution is still highly imperfect. On one hand, this does not guarantee the  
110 absence of fake accounts. On the other hand, and perhaps more importantly, this excludes most  
111 potential contributors from participating.

112 **Vouching mechanism.** To propagate trust to more accounts, Tournesol also proposes a vouching  
113 mechanism. Namely, any account can vouch for the authenticity of another account. More precisely,  
114 the account must vouch that the other account is used by a human who is not using any other account  
115 on the platform. The dataset contains 129 vouches.

116 **Comparison-based judgments.** Following a large literature on the topic [38, 17, 66, 10, 73, 60, 42],  
117 Tournesol relies on a comparison-based preference elicitation system. We believe that the need to  
118 distinguish among top content which should be more recommended makes this system more suitable  
119 than, e.g., using direct assessments [63, 2, 55, 85], which may yield too many “saturated” maximal  
120 assessments. Additionally, comparisons are labelled with the week in which the comparison was first  
121 submitted. This allows potentially observing changes or drifts in the contributors’ judgments.

122 Figure 1 (left) presents the video comparison interface. Namely, contributors are asked to select two  
123 videos, and to tell Tournesol which one of the videos should be recommended at scale. Moreover,  
124 rather than a binary decision, the contributor is asked to provide the judgment by moving a slider on  
125 a more continuous scale, from  $-10$  to  $10$ . The value  $-10$  means that the contributor would prefer  
126 Tournesol to recommend the left video vastly more often than the right videos, while the value  $0$   
127 means that they believe both videos should be recommended equally often.

128 **Quality criteria.** Tournesol allows contributors to rate nine other *optional* quality criteria (Figure 1)

- 129 • **Reliable and not misleading:** Is the presented information trustworthy, robustly backed and  
130 properly nuanced?
- 131 • **Clear and pedagogical:** How efficiently does the content guide viewers in their understanding?
- 132 • **Important and actionable:** Can additional focus on this topic have a significantly positive impact  
133 on the world?
- 134 • **Layman-friendly:** How understandable is it, without prior knowledge?

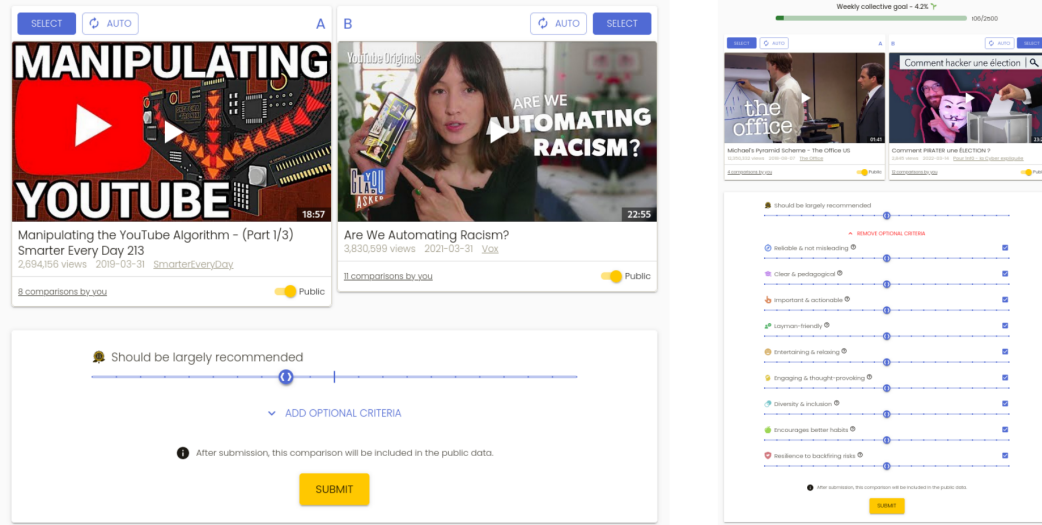


Figure 1: The interface through which contributors are asked to provide judgments. The judgments are comparisons of video contents using a slider along the main criteria "should be largely recommended" (left) and optional quality criteria (right).

- 135 • **Entertaining and relaxing:** Do people feel good watching it?
- 136 • **Engaging and thought-provoking:** Does it catch people's attention, spark curiosity and invite to
- 137 question previous beliefs?
- 138 • **Diversity and inclusion:** Does it promote tolerance, compassion and wider moral considerations?
- 139 • **Encourages better habits:** Does it make people adopt habits that benefit themselves and beyond?
- 140 • **Resilience to backfiring risks:** Is it adapted to viewers with opposing beliefs? Does it prevent
- 141 misconceptions or undesirable reactions?

142 While the criteria are further provided on Tournesol<sup>2</sup>, most contributors have surely *not* read thor-  
 143 oughly our descriptions. Arguably, they will more likely judge these criteria according to their own  
 144 understanding, which will be mostly based on the name of the criteria.

## 145 2.2 Processed data

146 In addition to the raw data presented thus far, the Tournesol public dataset exports processed data.  
 147 The processing is performed by a pipeline called SOLIDAGO [12].

148 **Solidago.** The pipeline has six modules. First, pretrust and vouches are used to assign *trust scores*  
 149 to all users. Second, *voting rights* are assigned to the different users, in a way that includes untrusted  
 150 users, while guaranteeing that they cannot outweigh trusted users. Third, for each criterion and each  
 151 user, the comparisons are turned into the user's *raw scores*, using the generalized Bradley-Terry  
 152 model [33]. Fourth, raw scores are *scaled*, using Mehestan [4], zero-shift and standardization. Fifth,  
 153 scaled scores are securely aggregated into *global scores*, using the Lipschitz-resilient quadratically  
 154 regularized quantile [12]. Sixth, all scores are squashed into  $(-100, 100)$ , using the map  $t \mapsto$   
 155  $100t/\sqrt{1+t^2}$ . All along, left and right uncertainties on all variables are computed.

156 **Exported values.** Trust scores, squashed individual scores and squashed global scores are provided  
 157 in the public dataset.

158 **Results.** Figure 2 lists the most recommendable videos, according to Tournesol's contributors, as  
 159 they are displayed on the website.

<sup>2</sup><https://tournesol.app/criteria>

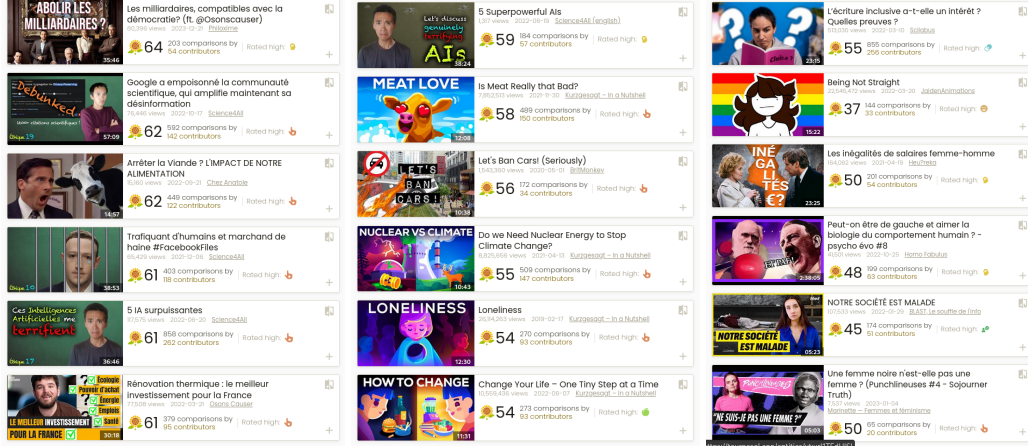


Figure 2: Best videos (left), best English-speaking videos (middle) and best videos along the criterion “diversity & inclusivity” (right).

## 2.3 Privacy

Overall, we encourage transparency in our contributors, as we believe that this will foster important research on human judgments, and help make safer and more ethical algorithms. However, we acknowledge that, because of social and political pressures, some judgments are dangerous to make public, e.g. when criticizing one’s own employer or government. This is why we allow contributors to provide data publicly or privately. More precisely, each contributor can select the privacy setting of any video they rate. If a video is rated privately, then all its comparisons to any other video will be recorded privately. Only Tournesol’s server can access to such data. Conversely, all comparisons that involve two publicly rated videos are exported in the Tournesol public dataset.

## 2.4 Data collection context

The contributors to Tournesol receive no financial compensation. Their contributions are mostly motivated by the desire to contribute to a democratic AI governance project, and by the will to promote content of public interest. Their recruitment is thus organic, and mostly depends on how frequently they were exposed to the promotion of the Tournesol project. Evidently, this greatly correlates with Tournesol’s communication, which has been heavily supported by the (French-speaking) YouTube channel Science4All, and by other science communicators [51]. As a result, the set of contributors is in no way representative of the global population. Namely, it is heavily biased towards science enthusiasts. Nevertheless, we believe that the data provided by this community should be of great interest to AI alignment, at least on topics with a significant scientific component.

## 3 Data analysis

This section presents some data analyses to provide insights in the *Tournesol public dataset*.

### 3.1 Contributors’ contributions

Figure 3 displays the number of contributions per user. Perhaps unsurprisingly, this statistics is heavy-tailed; in fact, it seems to fit Zipf’s law [82], with a few contributors providing most of the comparisons, and most of them providing very few. Figure 4 plots the activity through time: Tournesol has 100 to 200 weekly active users, while the number of monthly active users fluctuates between 200 and 900.

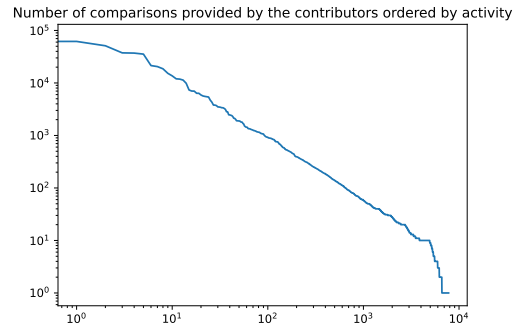


Figure 3: Number of comparisons provided by the different contributors, on a log-log scale, which is typical of Zipf’s law [82].

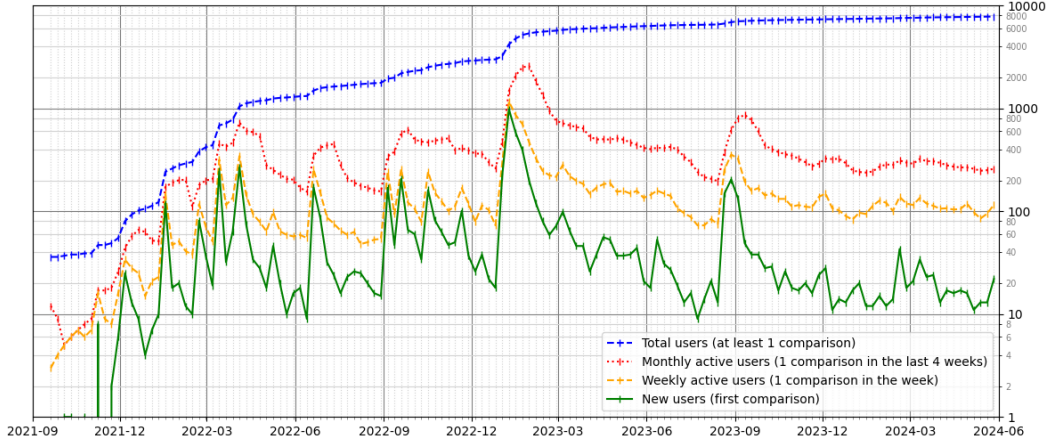


Figure 4: Contributors' participation through time.

### 3.2 Video and contributor connectivity

For scores to be meaningful, the contributors must have compared sufficiently many videos in common [4]. The contributor comparability graph has a connected component with 7187 contributors and diameter of 6, out of the 7,826 contributors that have compared at least 2 videos. The graph has 208,323 edges out of 30,619,225 possible (0.68%) making it very sparse. But for the induced graph of the top 100 most active contributors with a trust at least 0.1 (which correspond to *scaling-calibration* contributors [12]), 3,442 (69.5%) pairs of contributors are comparable. This justifies the restriction of scaling calibration to the most active contributors.

Figure 5 details video comparisons for some highly active users. Interestingly, because the platform lets contributors to select their videos to compare, we observe a wide variety of comparison graphs. This raises open questions about the uncertainties of the resulting learned scores [33], and about the possibility to improve accuracy through *active learning* [67, 83].

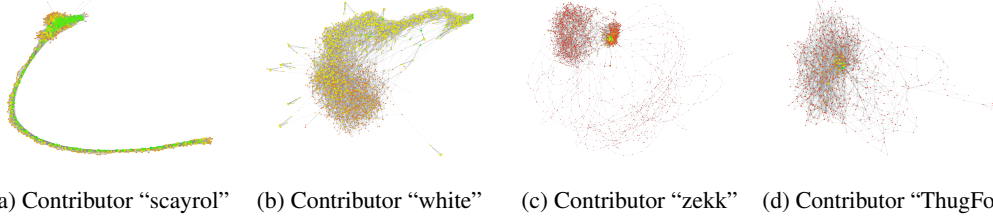


Figure 5: Graphs of video comparisons for different users

### 3.3 Correlations between criteria

Figure 6 reports the correlations between quality criteria, in contributors' comparative judgments. Perhaps most remarkably, we observe that the criterion that best predicts whether a video "should be more largely recommended" is whether it is "important and actionable". This finding highlights the need to pay greater attention to *information prioritization*, and especially combatting "*mute news*" [51]. In particular, there may be an excess of attention to "*fake news*". In fact, [75] expose numerous strategies from the "merchants of doubts" that do not involve producing false information, such as shifting blame, cyberbullying critics or "striking a positive tone" [24].

Figure 6 also shows that most criteria are only weakly correlated. Two notable exceptions are "important and actionable" and "encourage better habits", and "reliable and not misleading" and "clear and pedagogical", which could be argued to be slightly redundant.

Note also that, as expected given Berkson's paradox [11], the correlations decrease if we only consider the top 10% videos on Tournesol (i.e. those that are more likely to be recommended).

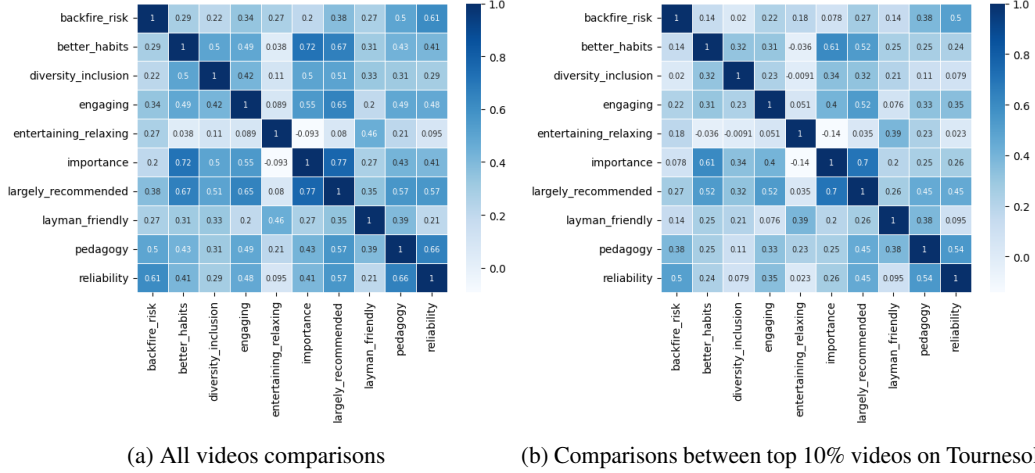


Figure 6: Correlations between quality criteria

### 3.4 Distributions of reported comparisons

As it is not formally defined how contributors should rate a pair of videos, we expected many different expression styles. We ran a clustering algorithm (K-means) on statistics of the distribution of comparison values for each user. Figure 7 shows the typical distribution of comparison values of each of the eight clusters we identified. While some contributors provided comparisons close to “recommend equally” (cluster 3 and 4), others’ comparisons were systematically towards the extreme (clusters 2, 5 and 6). This suggests that the discrepancies between their individual scores will be due to their expression style, rather than actual differences in their judgments, which justifies the research on mitigating the heterogeneity in expression styles [53, 93, 4].

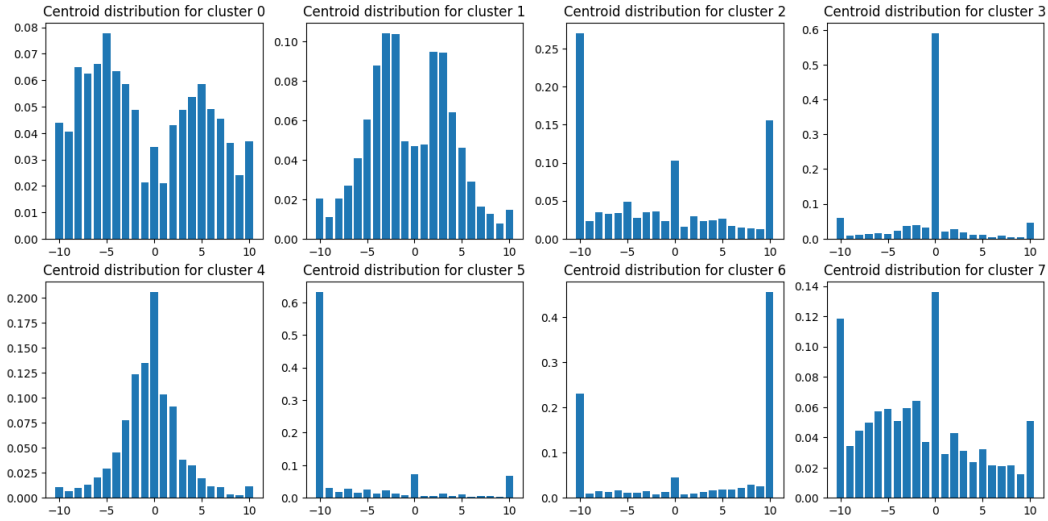


Figure 7: Example centroids of 8 clusters obtained by the K-means algorithm applied to the distributions of comparison values for each contributor with at least 20 comparisons.

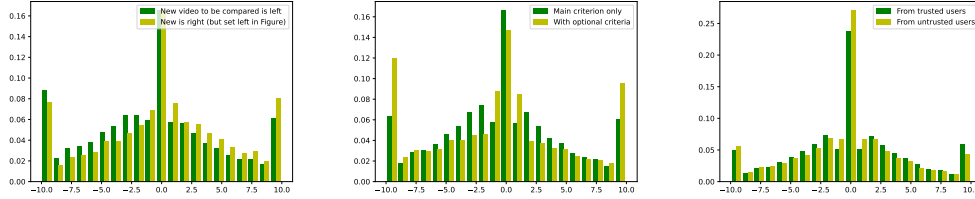
### 3.5 Psychological biases in contributors’ judgments

our dataset exposes psychological biases in contributors’ judgments. One example is a instinctive desire to over-recommend a recently watched high-quality video, known as the *recency bias* [62], which is depicted by Figure 8a. Namely, this figure plots all comparisons on the main criterion that correspond to a contributor evaluating a given video for the first time (negative scores correspond



to the newly scored videos). The 95% confidence interval for the mean of first-time comparisons is  $[-0.40, -0.32]$ , which is arguably a surprisingly significant bias.

Another bias we observe is a tendency to favor left videos. The 95% confidence interval for the mean of the main-criterion comparisons (Figure 8b) is  $[-0.49, -0.44]$ . Considering all criteria (Figure 8c) yields a smaller bias, with a corresponding 95% confidence interval of  $[-0.17, -0.15]$ . This suggests that reflecting on more criteria reduces the left-video bias. And indeed, when they are accompanied with comparisons on other criteria, the main-criterion comparisons have a 95% confidence interval for the mean equal to  $[-0.38, -0.31]$ , as opposed to  $[-0.57, 0.52]$  for main-criterion-only comparisons. We also observe that pretrusted contributors have a significantly reduced left-video bias (on all criteria,  $[-0.03, -0.002]$  for pretrusted,  $[-0.34, -0.31]$  for unpretrusted).



(a) First comparisons on main criterion, (b) Comparisons on main criterion, (c) Comparisons on all criteria, separated based on optional criteria, separated based on trust.

Figure 8: Recency and left-video biases in contributors' judgments.

### 3.6 Distribution of scores

Unsquashed scores (essentially, as outputs of the generalized Bradley-Terry model on contributors' comparisons) are extremely heavy tailed. Indeed, out of 634516 scores, 2803 deviate by more than 5 standard deviations. This is to be contrasted with the expected number 0.18 of such extreme scores, assuming a normal distribution of the scores. In fact, 428 scores deviate by more than 10 standard deviations. This observation justifies the use of comparisons to quantify the potential large deviations between top alternatives, which direct scoring approaches might fail to account for appropriately, as well as of a (robustified) quantile to standardize scores [12].

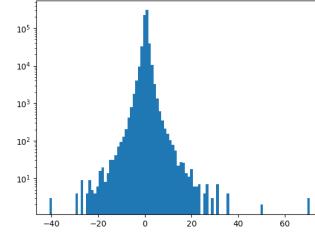


Figure 9: Distribution of unsquashed scores, with logarithmic y-scale.

## 4 Research challenges

Tournesol raises numerous fascinating research challenges. Below, we sketch some of these.

**Aggregate the different criteria into a score.** We expect the combination of many different quality criteria to yield a more reliable judgment of what content ought to be recommended at scale, or to a given specific user. However, the appropriate aggregation of our different quality criteria is still unclear, especially given probable nonlinear phenomena. How best to do this should be investigated.

**Debias the contributing population.** Like in many online participatory projects [9], we expect huge participation imbalances. Leveraging demographic data to debias the Tournesol recommendations, e.g., by giving more voting rights to individuals from underrepresented communities, could help, but it will require both (safely) collecting personal data and building new (secure) algorithms, akin to those used by the *Community Notes*<sup>3</sup> or by *Pol.is*<sup>4</sup>.

**Volition.** As Section 3.5 highlighted it, we cannot expect the Tournesol database to contain fully reliable human judgments. Many comparisons have surely been provided by contributors, at moments when they were not paying the utmost attention to all the possible ramifications and unwanted side

<sup>3</sup><https://communitynotes.twitter.com/guide/en/under-the-hood/ranking-notes>

<sup>4</sup><https://compdemocracy.org/algorithms/>



267 effects of promoting a video at scale. In particular, some judgments will arguably be more reliable  
268 than others. Such more reliable judgments are sometimes called *volitions*, rather than *preferences*.  
269 There is a need for algorithms that model human psychology to distinguish these two [50, 59].

270 **Privacy.** Tournesol’s current algorithms do not provide any *differential privacy* [26]. Future research  
271 should also investigate how to strengthen privacy without harming too much the quality and the  
272 security of the Tournesol scores. Perhaps most importantly, ideally, Tournesol’s servers would be  
273 able to leverage private comparisons to score videos without being a single point of failure for private  
274 data protection. *Secure multi-party computations* could be a promising venue to do so [20].

275 **Decentralize Tournesol.** A longer-term goal is to fully decentralize Tournesol. In this vision, the  
276 data would no longer be stored on Tournesol’s server, but would be replicated appropriately on a  
277 large number of contributors’ devices. Moreover, the computations of Tournesol scores should also  
278 be decentralized, while guaranteeing *Byzantine resilience* [58]. Recent research in fully decentralized  
279 Byzantine learning has provided the building blocks of such a decentralization [29, 35], but more  
280 research is needed to understand how to best do so in the context of Tournesol.

281 **Preference generalization.** Right now, contributors are only voting on the videos that they explicitly  
282 compared. However, if they consistently voted positively all the videos of a given channel, then we  
283 could guess that they would have voted positively a new video from this channel, and to include their  
284 likely vote even when they did not compare the new video. Evidently, additional information can  
285 be leveraged to make such generalizations, such as the other video features (description, transcript,  
286 length), and the other contributors’ judgments (using collaborative filtering [84]). Note however that  
287 generalization increases vulnerability risks. A careful security analysis would be required [68].

288 **Language model alignment.** Tournesol’s database could help align language models, e.g. through  
289 *reinforcement learning with Tournesol feedback* [21, 76]. Determining how to combine large language  
290 models [37] with Tournesol’s database to design safer models is an exciting venue for future work.

291 **Leverage expertise.** On technical topics like vaccination or climate change, especially when  
292 misconceptions are widespread in the general population, it seems desirable to assign more voting  
293 rights to experts, especially when judging the reliability of content within their domains of expertise.  
294 This issue is intimately connected to Condorcet’s jury problem [22, 72].

295 **Proof of Personhood with zero knowledge.** Combatting fake accounts arguably remains the top  
296 priority to secure participatory systems. To address this, at least in democratic countries and in the  
297 short term, the state could be tasked with delivering *Proofs of Personhood* [16, 41], if possible in a  
298 zero-knowledge manner. More precisely, any citizen should ideally be able to provide to any platform  
299 a proof of citizenship, which does not enable neither the platform nor the state to identify which  
300 account is owned by which citizen. We believe that designing such a system could have applications  
301 beyond the particular case of Tournesol. Indeed, we could demand that social media only display  
302 the number of likes from users with a delivered proof of citizenship, and that their recommendation  
303 algorithms be trained only by such certified users’ data.

304 **Liquid democracy** Finally, future work could investigate the extent to which a liquid democ-  
305 racy [48] could be set up on platforms like Tournesol. Such a system through which a contributor  
306 can delegate their votes to other voters could help combat activity bias (i.e. better accounting for  
307 inactive contributors) and expertise (if voters delegate to more competent contributors). While  
308 philosophically appealing, the security of such a system should however be first investigated [5].

## 309 5 Conclusion

310 This paper introduced the *Tournesol public dataset*, which is a large, secured and trustworthy database  
311 of reliable human judgments. We detailed its construction, and provided an analysis of its content.  
312 We believe that this database can help stimulate and facilitate research and development on ethical  
313 algorithms, and could eventually help improve the informational diet of billions of people for the better.  
314 Given the current information crisis, we regard this as an “important and actionable” contribution.

## References

- [1] Himan Abdollahpouri, Gediminas Adomavicius, Robin Burke, Ido Guy, Dietmar Jannach, Toshihiro Kamishima, Jan Krasnodebski, and Luiz Pizzato. Multistakeholder recommendation: Survey and research directions. *User Modeling and User-Adapted Interaction*, 30:127–158, 2020.
- [2] Gerald Albaum. The likert scale revisited. *Market Research Society. Journal.*, 39(2):1–21, 1997.
- [3] Richard P Allan, Paola A Arias, Sophie Berger, Josep G Canadell, Christophe Cassou, Deliang Chen, Annalisa Cherchi, Sarah L Connors, Erika Coppola, Faye Abigail Cruz, et al. Intergovernmental panel on climate change (ipcc). summary for policymakers. In *Climate change 2021: The physical science basis. Contribution of working group I to the sixth assessment report of the intergovernmental panel on climate change*, pages 3–32. Cambridge University Press, 2023.
- [4] Youssef Allouah, Rachid Guerraoui, Lê-Nguyên Hoang, and Oscar VILLEMAUD. Robust sparse voting. *CoRR*, abs/2202.08656, 2022.
- [5] Shiri Alouf-Heffetz, Tanmay Inamdar, Pallavi Jain, Nimrod Talmon, and Yash More Hiren. Controlling delegations in liquid democracy. In Mehdi Dastani, Jaime Simão Sichman, Natasha Alechina, and Virginia Dignum, editors, *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2024, Auckland, New Zealand, May 6-10, 2024*, pages 2624–2632. ACM, 2024.
- [6] Cécile Andrzejewski. “team jorge”: In the heart of a global disinformation machine. *Forbidden Stories*, 2023.
- [7] Fabio Angiolillo, Martin Lundstedt, Marina Nord, and Staffan I Lindberg. State of the world 2023: democracy winning and losing at the ballot. *Democratization*, pages 1–25, 2024.
- [8] Valentin Armhein, Sander Greenland, and Blake McShane. Scientists rise up against statistical significance. *Nature*, 567(7748):305–307, 2019.
- [9] Edmond Awad, Sohan Dsouza, Richard Kim, Jonathan Schulz, Joseph Henrich, Azim Shariff, Jean-François Bonnefon, and Iyad Rahwan. The moral machine experiment. *Nature*, 563(7729):59–64, 2018.
- [10] Edmond Awad, Sohan Dsouza, Richard Kim, Jonathan Schulz, Joseph Henrich, Azim Shariff, Jean-François Bonnefon, and Iyad Rahwan. The moral machine experiment. *Nature*, 563(7729):59–64, 2018.
- [11] Joseph Berkson. Limitations of the application of fourfold table analysis to hospital data. *Biometrics Bulletin*, 2(3):47–53, 1946.
- [12] Romain Beylerian, Bérangère Colbois, Louis Faucon, Lê Nguyễn Hoang, Aidan Jungo, Alain Le Noac’h, and Adrien Matissart. Tournesol: Permissionless collaborative algorithmic governance with security guarantees. *CoRR*, abs/2211.01179, 2022.
- [13] Peva Blanchard, El-Mahdi El-Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine learning with adversaries: Byzantine tolerant gradient descent. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 119–129, 2017.
- [14] Shawn Boburg. Leaked files reveal reputation-management firm’s deceptive tactics. *The Washington Post*, pages NA–NA, 2023.
- [15] Maria Borge, Eleftherios Kokoris-Kogias, Philipp Jovanovic, Linus Gasser, Nicolas Gailly, and Bryan Ford. Proof-of-personhood: Redemocratizing permissionless cryptocurrencies. In *2017 IEEE European Symposium on Security and Privacy Workshops, EuroS&P Workshops 2017, Paris, France, April 26-28, 2017*, pages 23–26. IEEE, 2017.

- [16] Maria Borge, Eleftherios Kokoris-Kogias, Philipp Jovanovic, Linus Gasser, Nicolas Gailly, and Bryan Ford. Proof-of-personhood: Redemocratizing permissionless cryptocurrencies. In *2017 IEEE European Symposium on Security and Privacy Workshops, EuroS&P Workshops 2017, Paris, France, April 26-28, 2017*, pages 23–26. IEEE, 2017.
- [17] Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- [18] Samantha Bradshaw and Philip N Howard. The global organization of social media disinformation campaigns. *Journal of International Affairs*, 71(1.5):23–32, 2018.
- [19] Luca Braghieri, Ro’ee Levy, and Alexey Makarin. Social media and mental health. *American Economic Review*, 112(11):3660–3693, 2022.
- [20] Ran Canetti, Uriel Feige, Oded Goldreich, and Moni Naor. Adaptively secure multi-party computation. In Gary L. Miller, editor, *Proceedings of the Twenty-Eighth Annual ACM Symposium on the Theory of Computing, Philadelphia, Pennsylvania, USA, May 22-24, 1996*, pages 639–648. ACM, 1996.
- [21] Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4299–4307, 2017.
- [22] Marie Jean Antoine Nicolas de Caritat Condorcet. *Essai sur l’application de l’analyse à la probabilité des décisions rendues à la pluralité des voix*. L’imprimerie royale, 1785.
- [23] Paul Covington, Jay Adams, and Emre Sargin. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM conference on recommender systems*, pages 191–198, 2016.
- [24] Paresh Dave and Jeffrey Dastin. Google told its scientists to ‘strike a positive tone’ in ai research - documents. *Reuters*, 2023.
- [25] John R Douceur. The sybil attack. In *International workshop on peer-to-peer systems*, pages 251–260. Springer, 2002.
- [26] Cynthia Dwork. Differential privacy. In Michele Bugliesi, Bart Preneel, Vladimiro Sassone, and Ingo Wegener, editors, *Automata, Languages and Programming, 33rd International Colloquium, ICALP 2006, Venice, Italy, July 10-14, 2006, Proceedings, Part II*, volume 4052 of *Lecture Notes in Computer Science*, pages 1–12. Springer, 2006.
- [27] El Mahdi El Mhamdi. *Robust Distributed Learning*. PhD thesis, EPFL, 2020.
- [28] El-Mahdi El-Mhamdi, Sadegh Farhadkhani, Rachid Guerraoui, Arsany Guirguis, Lê Nguyễn Hoang, and Sébastien Rouault. Collaborative learning as an agreement problem. *CoRR*, abs/2008.00742, 2020.
- [29] El-Mahdi El-Mhamdi, Sadegh Farhadkhani, Rachid Guerraoui, Arsany Guirguis, Lê Nguyễn Hoang, and Sébastien Rouault. Collaborative learning in the jungle. *CoRR*, abs/2008.00742, 2020.
- [30] El-Mahdi El-Mhamdi, Sadegh Farhadkhani, Rachid Guerraoui, Nirupam Gupta, Lê-Nguyễn Hoang, Rafael Pinot, and John Stephan. On the impossible safety of large AI models. *CoRR*, abs/2209.15259, 2022.
- [31] El-Mahdi El-Mhamdi, Rachid Guerraoui, and Sébastien Rouault. The hidden vulnerability of distributed learning in byzantium. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 3518–3527. PMLR, 2018.

- 410 [32] El-Mahdi El-Mhamdi, Rachid Guerraoui, and Sébastien Rouault. Distributed momentum for  
411 byzantine-resilient learning. *CoRR*, abs/2003.00010, 2020.
- 412 [33] Julien Fageot, Sadegh Farhadkhani, Lê-Nguyễn Hoang, and Oscar VILLEMAUD. Generalized  
413 bradley-terry models for score estimation from paired comparisons. In Michael J. Wooldridge,  
414 Jennifer G. Dy, and Sriraam Natarajan, editors, *Thirty-Eighth AAAI Conference on Artificial  
415 Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelli-  
416 gence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence,  
417 EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 20379–20386. AAAI Press, 2024.
- 418 [34] Jack Farchy. Oil trader sues uae claiming smear campaign bankrupted his firm. *Bloomberg*,  
419 2024.
- 420 [35] Sadegh Farhadkhani, Rachid Guerraoui, Nirupam Gupta, Lê-Nguyễn Hoang, Rafael Pinot,  
421 and John Stephan. Robust collaborative learning with linear gradient overhead. In Andreas  
422 Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan  
423 Scarlett, editors, *International Conference on Machine Learning, ICML 2023, 23-29 July 2023,  
424 Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages  
425 9761–9813. PMLR, 2023.
- 426 [36] Sadegh Farhadkhani, Rachid Guerraoui, Lê Nguyễn Hoang, and Oscar VILLEMAUD. An equiva-  
427 lence between data poisoning and byzantine gradient attacks. In Kamalika Chaudhuri, Stefanie  
428 Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International Confer-  
429 ence on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume  
430 162 of *Proceedings of Machine Learning Research*, pages 6284–6323. PMLR, 2022.
- 431 [37] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion  
432 parameter models with simple and efficient sparsity. *CoRR*, abs/2101.03961, 2021.
- 433 [38] Leon Festinger. A theory of social comparison processes. *Human relations*, 7(2):117–140,  
434 1954.
- 435 [39] Christina Fink. Dangerous speech, anti-muslim violence, and facebook in myanmar. *Journal of  
436 International Affairs*, 71(1.5):43–52, 2018.
- 437 [40] Philippa Foot. The problem of abortion and the doctrine of double effect. *Oxford Review*, 5,  
438 1967.
- 439 [41] Bryan Ford. Identity and personhood in digital democracy: Evaluating inclusion, equality, secu-  
440 rity, and privacy in pseudonym parties and other proofs of personhood. *CoRR*, abs/2011.02412,  
441 2020.
- 442 [42] Rachel Freedman, Jana Schaich Borg, Walter Sinnott-Armstrong, John P. Dickerson, and  
443 Vincent Conitzer. Adapting a kidney exchange algorithm to align with human values. *Artif.  
444 Intell.*, 283:103261, 2020.
- 445 [43] Elia Gabarron, Sunday Oluwafemi Oyeyemi, and Rolf Wynn. Covid-19-related misinformation  
446 on social media: a systematic review. *Bulletin of the World Health Organization*, 99(6):455,  
447 2021.
- 448 [44] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna M.  
449 Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets. *Commun. ACM*, 64(12):86–  
450 92, 2021.
- 451 [45] Dominique Geissler, Dominik Bär, Nicolas Pröllochs, and Stefan Feuerriegel. Russian propa-  
452 ganda on social media during the 2022 invasion of ukraine. *EPJ Data Science*, 12(1):35,  
453 2023.
- 454 [46] Shane Griffith, Kaushik Subramanian, Jonathan Scholz, Charles L. Isbell Jr., and Andrea Lock-  
455 erd Thomaz. Policy shaping: Integrating human feedback with reinforcement learning. In  
456 Christopher J. C. Burges, Léon Bottou, Zoubin Ghahramani, and Kilian Q. Weinberger, editors,  
457 *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural  
458 Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake  
459 Tahoe, Nevada, United States*, pages 2625–2633, 2013.

- [47] Gillian Kereldena Hadfield. *Rules for a flat world: why humans invented law and how to reinvent it for a complex global economy*. Oxford University Press, 2017.
- [48] Daniel Halpern, Joseph Y. Halpern, Ali Jadbabaie, Elchanan Mossel, Ariel D. Procaccia, and Manon Revel. In defense of liquid democracy. In Kevin Leyton-Brown, Jason D. Hartline, and Larry Samuelson, editors, *Proceedings of the 24th ACM Conference on Economics and Computation, EC 2023, London, United Kingdom, July 9-12, 2023*, page 852. ACM, 2023.
- [49] Adam D Hernandez. Cambridge analytica. *Class, Race and Corporate Power*, 11(2), 2023.
- [50] Lê Nguyễn Hoang. Towards robust end-to-end alignment. In Huáscar Espinoza, Seán Ó hÉigeartaigh, Xiaowei Huang, José Hernández-Orallo, and Mauricio Castillo-Effen, editors, *Workshop on Artificial Intelligence Safety 2019 co-located with the Thirty-Third AAAI Conference on Artificial Intelligence 2019 (AAAI-19), Honolulu, Hawaii, January 27, 2019*, volume 2301 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2019.
- [51] Lê Nguyễn Hoang. Science communication desperately needs more aligned recommendation algorithms. *Frontiers in Communication*, 5:115, 2020.
- [52] Le Nguyen Hoang and El Mahdi El Mhamdi. *Le fabuleux chantier: Rendre l'intelligence artificielle robustement bénéfique*. edp Sciences, 2019.
- [53] Lê Nguyễn Hoang, François Soumis, and Georges Zaccour. Measuring unfairness feeling in allocation problems. *Omega*, 65:138–147, 2016.
- [54] Chiungjung Huang. A meta-analysis of the problematic social media use and mental health. *International Journal of Social Psychiatry*, 68(1):12–33, 2022.
- [55] Ankur Joshi, Saket Kale, Satish Chandel, and D Kumar Pal. Likert scale: Explored and explained. *Current Journal of Applied Science and Technology*, pages 396–403, 2015.
- [56] Jastra Kanjec. Facebook removed more than 15 billion fake accounts in two years, five times more than its active user base. *StockApps*, 2021.
- [57] Sai Praneeth Karimireddy, Lie He, and Martin Jaggi. Byzantine-robust learning on heterogeneous datasets via bucketing. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.
- [58] Leslie Lamport, Robert E. Shostak, and Marshall C. Pease. The byzantine generals problem. *ACM Trans. Program. Lang. Syst.*, 4(3):382–401, 1982.
- [59] Mohamed Lechiakh and Alexandre Maurer. Volition learning: What would you prefer to prefer? In Helmut Degen and Stavroula Ntoa, editors, *Artificial Intelligence in HCI - 4th International Conference, AI-HCI 2023, Held as Part of the 25th HCI International Conference, HCII 2023, Copenhagen, Denmark, July 23-28, 2023, Proceedings, Part I*, volume 14050 of *Lecture Notes in Computer Science*, pages 555–574. Springer, 2023.
- [60] Min Kyung Lee, Daniel Kusbit, Anson Kahng, Ji Tae Kim, Xinran Yuan, Allissa Chan, Daniel See, Ritesh Noothigattu, Siheon Lee, Alexandros Psomas, and Ariel D. Procaccia. Webuidai: Participatory framework for algorithmic governance. *Proc. ACM Hum. Comput. Interact.*, 3(CSCW):181:1–181:35, 2019.
- [61] Xiangru Lian, Binhang Yuan, Xuefeng Zhu, Yulong Wang, Yongjun He, Honghuan Wu, Lei Sun, Haodong Lyu, Chengjun Liu, Xing Dong, et al. Persia: An open, hybrid system scaling deep learning-based recommenders up to 100 trillion parameters. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 3288–3298, 2022.
- [62] David A Liebermann. *Learning and memory: An integrative approach*. Belmont, CA: Thomson/Wadsworth, 2004.
- [63] Rensis Likert. A technique for the measurement of attitudes. *Archives of psychology*, 1932.
- [64] Zhifan Luo. “why should facebook (not) ban trump?”: connecting divides in reasoning and morality in public deliberation. *Information, Communication & Society*, 25(5):654–668, 2022.

- [65] Kirsten Martin. Recommending an insurrection: Facebook and recommendation algorithms. In *Ethics of Data and Analytics*, pages 225–239. Auerbach Publications, 2022.
- [66] Lucas Maystre. *Efficient Learning from Comparisons*. PhD thesis, EPFL, 2018.
- [67] Lucas Maystre and Matthias Grossglauser. Just sort it! A simple and effective approach to active preference learning. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 2344–2353. PMLR, 2017.
- [68] Bhaskar Mehta and Thomas Hofmann. A survey of attack-resistant collaborative filtering algorithms. *IEEE Data Eng. Bull.*, 31(2):14–22, 2008.
- [69] Silvia Milano, Mariarosaria Taddeo, and Luciano Floridi. Ethical aspects of multi-stakeholder recommendation systems. *The information society*, 37(1):35–45, 2021.
- [70] Michael K Miller. A republic, if you can keep it: Breakdown and erosion in modern democracies. *The Journal of Politics*, 83(1):198–213, 2021.
- [71] Paul Mozur. A genocide incited on facebook, with posts from myanmar’s military. *The New York Times*, 15(10):2018, 2018.
- [72] Shmuel Nitzan and Jacob Paroush. Optimal decision rules in uncertain dichotomous choice situations. *International Economic Review*, pages 289–297, 1982.
- [73] Ritesh Noothigattu, Snehal Kumar (Neil) S. Gaikwad, Edmond Awad, Sohan Dsouza, Iyad Rahwan, Pradeep Ravikumar, and Ariel D. Procaccia. A voting-based system for ethical decision making. In Sheila A. McIlraith and Kilian Q. Weinberger, editors, *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 1587–1594. AAAI Press, 2018.
- [74] Alina Oprea and Apostol Vassilev. Adversarial machine learning: A taxonomy and terminology of attacks and mitigations. Technical report, National Institute of Standards and Technology, 2023.
- [75] Naomi Oreskes and Erik M Conway. *Merchants of doubt: How a handful of scientists obscured the truth on issues from tobacco smoke to global warming*. Bloomsbury Publishing USA, 2011.
- [76] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [77] Luca Righes, Mohammed Saeed, Gianluca Demartini, and Paolo Papotti. The community notes observatory: Can crowdsourced fact-checking be trusted in practice? In Ying Ding, Jie Tang, Juan F. Sequeda, Lora Aroyo, Carlos Castillo, and Geert-Jan Houben, editors, *Companion Proceedings of the ACM Web Conference 2023, WWW 2023, Austin, TX, USA, 30 April 2023 - 4 May 2023*, pages 172–175. ACM, 2023.
- [78] Yasmim Mendes Rocha, Gabriel Acácio de Moura, Gabriel Alves Desidério, Carlos Henrique de Oliveira, Francisco Dantas Lourenço, and Larissa Deadame de Figueiredo Nicolete. The impact of fake news on social media and its influence on health during the covid-19 pandemic: A systematic review. *Journal of Public Health*, pages 1–10, 2021.
- [79] Allen L. Schirm Ronald Wasserstein and Nicole A. Lazar. Moving to a world beyond “ $p < 0.05$ ”. *The American Statistician*, 73:1–19, 2019.
- [80] Sébastien Rouault. *Practical Byzantine-resilient Stochastic Gradient Descent*. PhD thesis, EPFL, 2021.

- 556 [81] Stuart Russell. *Human compatible: Artificial intelligence and the problem of control*. Penguin,  
557 2019.
- 558 [82] Alexander I. Saichev, Yannick Malevergne, and Didier Sornette. *Theory of Zipf’s law and  
559 beyond*, volume 632. Springer Science & Business Media, 2009.
- 560 [83] Ayush Sekhari, Karthik Sridharan, Wen Sun, and Runzhe Wu. Contextual bandits and imitation  
561 learning with preference-based active queries. In Alice Oh, Tristan Naumann, Amir Globerson,  
562 Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information  
563 Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023,  
564 NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- 565 [84] Xiaoyuan Su and Taghi M. Khoshgoftaar. A survey of collaborative filtering techniques. *Adv.  
566 Artif. Intell.*, 2009:421425:1–421425:19, 2009.
- 567 [85] Basu Prasad Subedi. Using likert type data in social science research: Confusion, issues and  
568 challenges. *International journal of contemporary applied sciences*, 3(2):36–49, 2016.
- 569 [86] Gan Sun, Yang Cong, Jiahua Dong, Qiang Wang, Lingjuan Lyu, and Ji Liu. Data poisoning  
570 attacks on federated machine learning. *IEEE Internet of Things Journal*, 9(13):11365–11375,  
571 2021.
- 572 [87] Maxime Tellier. Enquête avisa partners : dans les coulisses de la sulfureuse agence d’influence  
573 soupçonnée de désinformation. *France Info*, 2023.
- 574 [88] Mariame Tighanimine. *L’affaiblissement des corps intermédiaires par les plateformes Internet.  
575 Le cas des médias et des syndicats français au moment des Gilets jaunes*. Conservatoire National  
576 des Arts et Métiers, 2019.
- 577 [89] Petter Törnberg. How digital media drive affective polarization through partisan sorting.  
578 *Proceedings of the National Academy of Sciences*, 119(42):e2207159119, 2022.
- 579 [90] Zeynep Tufekci. *Twitter and tear gas: The power and fragility of networked protest*. Yale  
580 University Press, 2017.
- 581 [91] Jean M Twenge, Jonathan Haidt, Jimmy Lozano, and Kevin M Cummins. Specification curve  
582 analysis shows that social media use is linked to poor mental health, especially among girls.  
583 *Acta psychologica*, 224:103512, 2022.
- 584 [92] Myriam Vidal Valero. Thousands of scientists are cutting back on twitter. *Nature*, 620:482–4,  
585 2023.
- 586 [93] Jingyan Wang and Nihar B. Shah. Your 2 is my 1, your 3 is my 9: Handling arbitrary  
587 miscalibrations in ratings. In Edith Elkind, Manuela Veloso, Noa Agmon, and Matthew E.  
588 Taylor, editors, *Proceedings of the 18th International Conference on Autonomous Agents and  
589 MultiAgent Systems, AAMAS ’19, Montreal, QC, Canada, May 13-17, 2019*, pages 864–872.  
590 International Foundation for Autonomous Agents and Multiagent Systems, 2019.
- 591 [94] Gabriel Weimann and Natalie Masri. Research note: Spreading hate on tiktok. *Studies in  
592 conflict & terrorism*, 46(5):752–765, 2023.
- 593 [95] Valerie Wirtschafter and Sharanya Majumder. Future challenges for online, crowdsourced  
594 content moderation: Evidence from twitter’s community notes. *Journal of Online Trust and  
595 Safety*, 2(1), Sep. 2023.



## 596 A Datasheet for the Tournesol dataset

597 In this appendix, we provide a datasheet for the Tournesol dataset, based on the framework proposed  
598 by [44].

### 599 A.1 Motivation

600 **For what purpose was the dataset created?** The dataset was created to identify videos of public  
601 interest that should be recommended more largely. Additionally, we hope that the dataset will help  
602 motivate research on the ethics and security of recommendation algorithms.

603 **Who created the dataset and on behalf of which entity?** The dataset was created by the nonprofit  
604 Tournesol Association, which is based in Switzerland.

605 **Who funded the creation of the dataset?** The Tournesol Association is supporting the creation  
606 and maintenance of the dataset. It is in majority funded by crowdsourced donations, with occasional  
607 services to private companies.

### 608 A.2 Composition

609 **What do the instances that comprise the dataset represent?** The dataset contains mostly pairwise  
610 comparisons of videos by users. The dataset also contains vouches between users, authentication  
611 status, as well as processed data from this raw data.

612 **How many instances are there in total?** The dataset contains 20k users (703 pretrusted), 40k  
613 videos, 126 vouches, 204k comparisons along the main criterion and 703k comparisons along optional  
614 criteria.

615 **Does the dataset contain all possible instances or is it a sample of instances of a larger set?** The  
616 dataset contains all *public* judgments provided on the Tournesol platform.

617 **What data does each instance consist of?** Each user has a pretrust status, based on email domain  
618 Sybil resilience. Each comparison is along a criterion, and refers to a user and a pair of videos.

619 **Is there a label or target associated with each instance?** Each comparison takes a value between  
620 -10 and 10.

621 **Is any information missing from individual instances?** Yes, plenty, such as the time it took to  
622 provide an answer, whether it was provided on a phone or a desktop, or whether the contributor  
623 actually watched the compared videos.

624 **Are relationships between individual instances made explicit?** Some of them, yes, such as the  
625 contributor’s identifier, or the videos that are compared.

626 **Are there recommended data splits?** Yes, comparisons are naturally split by criterion, or by users.  
627 Trusted/untrusted contributions could be split.

628 **Are there any errors, sources of noise, or redundancies in the dataset?** The comparisons come  
629 from humans, and are thus noisy, as well as potentially biased as discussed in the main part of the  
630 paper. Note that 4,446 comparisons were made before January 11, 2021, but because of a migration  
631 of the code, are dated on the January 11, 2021 week.

632 **Is the dataset self-contained, or does it link to or otherwise rely on external sources?** The  
633 dataset refers to YouTube videos, but could be analyzed without knowledge of the videos.

634 **Does the dataset contain data that might be considered confidential?** No. It was designed to be  
635 public.

636 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening**  
637 **or might otherwise cause anxiety?** Some poorly scored videos could be of this sort. Their content  
638 is not directly in the dataset, but the dataset points to them.

639 **Does the dataset identify any subpopulations?** Yes, trusted and untrusted contributors.

640 **Is it possible to identify individuals, either directly or indirectly, from the dataset?** Yes,  
641 especially given their public usernames.

642 **Does the dataset contain data that might be considered sensitive in any way?** Yes, indirectly, as  
643 it reveals consumption habits of contributors.

644 **Any other comments?** The individuals not only gave their consent, but the Tournesol also aims to  
645 make it clear that their provided data are used to design a democratic governance, and as such, could  
646 and should be scrutinized.

### 647 **A.3 Collection process**

648 **How was the data associated with each instance acquired?** Through the Tournesol platform  
649 <https://tournesol.app>.

650 **What mechanisms or procedures were used to collect the data?** Through the Tournesol compar-  
651 ison interface <https://tournesol.app/comparison>.

652 **If the dataset is a sample from a larger set, what was the sampling strategy?** Based on  
653 public/private settings selected by the contributor.

654 **Who was involved in the data collection process and how were they compensated?** Contributors  
655 are volunteers, most of whom are recruited through promotion in science YouTube videos. They are  
656 not compensated.

657 **Over what timeframe was the data collected?** The first data was collected in May 2020. The  
658 collection has been continuously ongoing since.

659 **Were any ethical review processes conducted?** Not by an institutional review board, as our work  
660 was done by a nonprofit association.

661 **Did you collect the data from the individuals in question directly, or obtain it via third parties**  
662 **or other sources?** Yes, through the Tournesol platform that we designed.

663 **Were the individuals in question notified about the data collection?** Yes. They had to cre-  
664 ate a Tournesol account, to consent with the data collection, and to select whether to make their  
665 contributions public or not.

666 **Did the individuals in question consent to the collection and use of their data?** Yes.

667 **If consent was obtained, were the consenting individuals provided with a mechanism to revoke**  
668 **their consent in the future or for certain uses?** Yes, contributors can delete their Tournesol  
669 account, which will delete their data from Tournesol’s (public) dataset.

670 **Has an analysis of the potential impact of the dataset and its use on data subjects been con-**  
671 **ducted?** Yes, we are consistently trying to make our project robustly beneficial.

### 672 **A.4 Preprocessing/cleaning/labeling**

673 **Was any preprocessing/cleaning/labeling of the data done?** Yes. To output trust scores, as well  
674 as squashed individual and global scores.

675 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data?** Yes. It is  
676 published in the Tournesol dataset.

677 **Is the software that was used to preprocess/clean/label the data available?** Yes. It is the  
678 open-source free-license Solidago python package.

## 679 **A.5 Uses**

680 **Has the dataset been used for any tasks already?** Yes, it is used to make content recommendations  
681 to 10k+ users.

682 **Is there a repository that links to any or all papers or systems that use the dataset?** Such  
683 papers and systems are listed in `tournesol.app/#research`.

684 **What (other) tasks could the dataset be used for?**

685 **Is there anything about the composition of the dataset or the way it was collected and prepro-**  
686 **cessed/cleaned/labeled that might impact future uses?**

687 **Are there tasks for which the dataset should not be used?** The dataset should not be used to  
688 harm individuals, communities or society.

## 689 **A.6 Distribution**

690 **Will the dataset be distributed to third parties outside of the entity (e.g., company, insti-**  
691 **tution, organization) on behalf of which the dataset was created?** Yes. It is published on  
692 `api.tournesol.app/exports/all`.

693 **How will the dataset be distributed?** zip file downloadable from the website.

694 **When will the dataset be distributed?** Already is.

695 **Will the dataset be distributed under a copyright or other intellectual property license, and/or**  
696 **under applicable terms of use?** Yes, it is under ODC-By license.

697 **Have any third parties imposed IP-based or other restrictions on the data associated with the**  
698 **instances?** No.

699 **Do any export controls or other regulatory restrictions apply to the dataset or to individual**  
700 **instances?** Not to our knowledge.

## 701 **A.7 Maintenance**

702 **Who will be supporting/hosting/maintaining the dataset?** The Tournesol association.

703 **How can the owner/curator/manager of the dataset be contacted?** `hello@tournesol.app`

704 **Is there an erratum?** No.

705 **Will the dataset be updated?** Yes. It is weekly updated, based on Tournesol’s users newly reported  
706 data.

707 **If the dataset relates to people, are there applicable limits on the retention of the data associated**  
708 **with the instances?** No limit applies.

709 **Will older versions of the dataset continue to be supported/hosted/maintained?** Yes, the dataset  
710 is consistently updated every week, based on contributors’ activity.

711 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**  
712 **them to do so?** The dataset is fully under the control of the Tournesol association. It is however  
713 under ODC-By license, thus any reuse is welcome, as long as attribution is appropriately provided.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [Yes]

Justification: The main contribution is, as explained, the publication of the dataset.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We explained the context in which the data is provided, and the limitations that this implies.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: Our paper does not provide theoretical results.

### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The code base and the data is available online and under copyleft free license.

### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The data is available at <https://api.tournesol.app/exports/all>, and the code is available at <https://github.com/tournesol-app/tournesol/>.

### 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We

### 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We did not provide statistical significance measures, mostly because statistical significance has been heavily criticized [79, 8]. Instead, we reported 95% confidence intervals. Note that the fact that they do not contain some “null hypothesis” is equivalent to saying that the null hypothesis has an associated p-value less than 5%. However, we believe that reporting confidence intervals is more meaningful, as it also communicates the effect size and an estimate of the uncertainty on the effect size.

### 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

763 Answer: [No]

764 Justification: No significant compute resource is needed. The graphs were all produced on

765 basic machines, without the need of, e.g., a GPU.

766 **9. Code Of Ethics**

767 Question: Does the research conducted in the paper conform, in every respect, with the

768 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

769 Answer: [Yes]

770 Justification: Our data collection platform <https://tournesol.app> repeatedly stresses

771 the fact that it aims to collect a public dataset of human judgments to help research. Explicit

772 consent is asked when contributors create their account. We make it clear that the contri-

773 butions should be made on a voluntarily basis, to help improve the security and ethics of

774 recommendation algorithms.

775 **10. Broader Impacts**

776 Question: Does the paper discuss both potential positive societal impacts and negative

777 societal impacts of the work performed?

778 Answer: [Yes]

779 Justification: The Tournesol project is fully motivated by the desire to have a positive societal

780 impact, by advancing the frontier of the research on the governance of recommendation

781 algorithms. We believe that these positive impacts clearly outweigh, and by far, the potential

782 negative societal impact, which could include, for instance, the ability of cybercrime to

783 better organize themselves.

784 **11. Safeguards**

785 Question: Does the paper describe safeguards that have been put in place for responsible

786 release of data or models that have a high risk for misuse (e.g., pretrained language models,

787 image generators, or scraped datasets)?

788 Answer: [Yes]

789 Justification: The dataset carefully annotates the source of the data, and contains information

790 on the degree of authentication of the sources.

791 **12. Licenses for existing assets**

792 Question: Are the creators or original owners of assets (e.g., code, data, models), used in

793 the paper, properly credited and are the license and terms of use explicitly mentioned and

794 properly respected?

795 Answer: [Yes]

796 Justification: The dataset is published by ourselves, under ODC-By license.

797 **13. New Assets**

798 Question: Are new assets introduced in the paper well documented and is the documentation

799 provided alongside the assets?

800 Answer: [Yes]

801 Justification: The dataset is documented in the paper, and a datasheet for datasets is provided

802 in the appendix.

803 **14. Crowdsourcing and Research with Human Subjects**

804 Question: For crowdsourcing experiments and research with human subjects, does the paper

805 include the full text of instructions given to participants and screenshots, if applicable, as

806 well as details about compensation (if any)?

807 Answer: [Yes]

808 Justification: We provided screenshots and contextualized the data collection process.

809 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**

810 **Subjects**

811 Question: Does the paper describe potential risks incurred by study participants, whether  
812 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
813 approvals (or an equivalent approval/review based on the requirements of your country or  
814 institution) were obtained?

815 Answer: [\[Yes\]](#)

816 Justification: The research was conducted by a nonprofit Association, and did not involve an  
817 IRB. We discussed the main risk for participants, namely retaliation from the entities they  
818 criticize. We stress, however, that this is usually not increasing the risk, compared to what  
819 they may already be publishing on social media.