

# Secure statistics for collaborative algorithmic governance

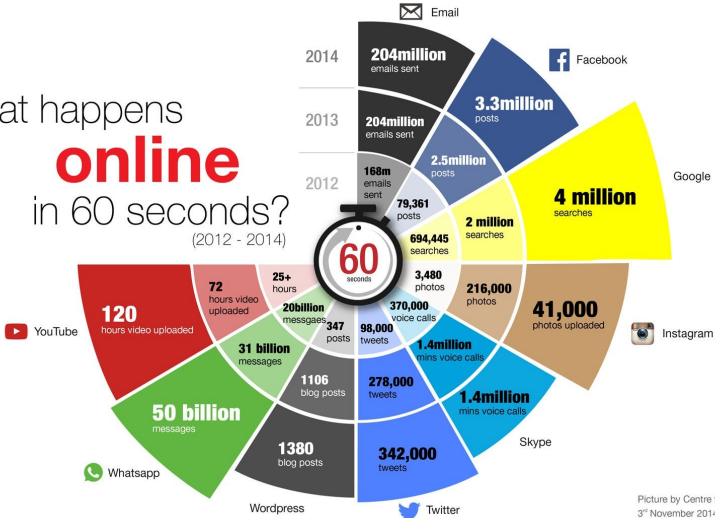
Lê Nguyễn Hoàng,  
Calicarpa, Tournesol & Science4All,  
STATLEARN, April 2023



## Section 1

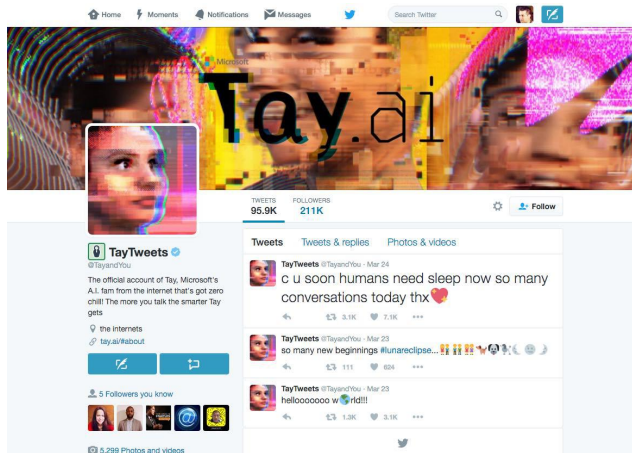
# Adversarial statistics

What happens  
**online**  
in 60 seconds?  
(2012 - 2014)



Picture by Centre for Learning and Teaching  
3<sup>rd</sup> November 2014

# The tale of Microsoft's two sisters (Tay vs Xiaoice)



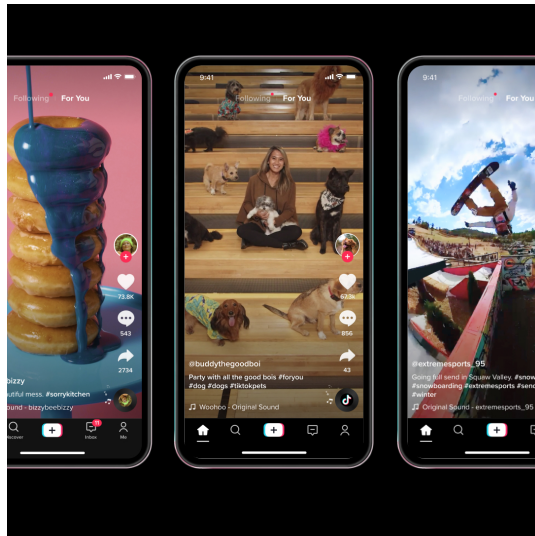
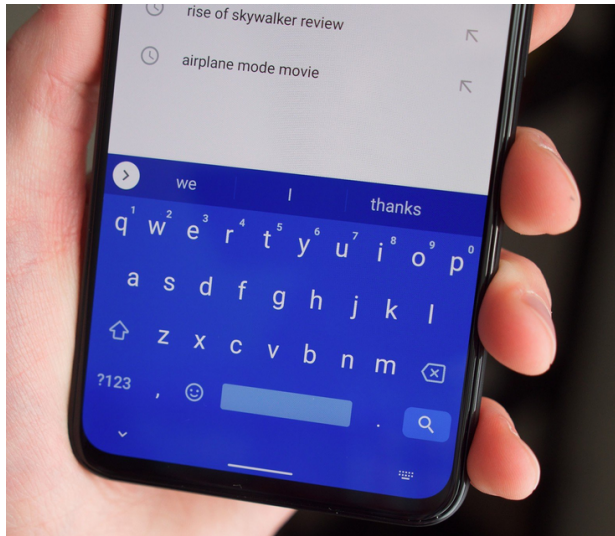
The screenshot shows the Twitter profile for Tay (@TayandYou). The header features a vibrant, abstract background with the text "Tay.ai" in a stylized font. The profile name is "TayTweets" with a verified badge, and the handle is "@TayandYou". The bio reads: "The official account of Tay, Microsoft's A.I. fam from the internet that's got zero chill! The more you talk the smarter Tay gets". The profile has 95.9K tweets and 211K followers. Three tweets are visible:

- Tweet 1: "c u soon humans need sleep now so many conversations today thx" (3.1K retweets, 7.1K likes)
- Tweet 2: "so many new beginnings #lunareclipse..." (111 retweets, 604 likes)
- Tweet 3: "helloooooo world!!!" (1.3K retweets, 3.1K likes)

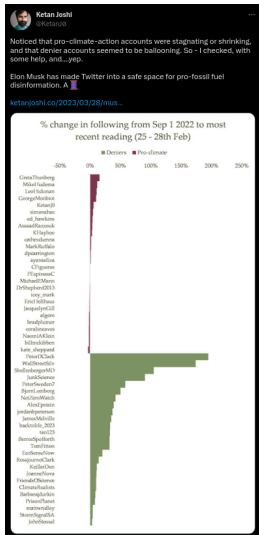




# You are voting all the times



# Non-users are stakeholders in online votes



AMNESTY INTERNATIONAL ENGLISH

WHO WE ARE WHAT WE DO COUNTRIES GET INVOLVED LATEST DONATE NOW

## MYANMAR: FACEBOOK'S SYSTEMS PROMOTED VIOLENCE AGAINST ROHINGYA; META OWES REPARATIONS

ACT NOW

News MYANMAR PRESS RELEASE SOUTH EAST ASIA AND THE PACIFIC TECHNOLOGY AND HUMAN RIGHTS

© Amnesty Intern

September 29, 2022

Facebook owner Meta's dangerous algorithms and reckless pursuit of profit substantially contributed to the atrocities perpetrated by the Myanmar military against the Rohingya people in 2017, [Amnesty International said in a new report published today](#).



# There is a huge steal-the-online-vote industry

## “TEAM JORGE”: IN THE HEART OF A GLOBAL DISINFORMATION MACHINE

In Part 2 of the “Story Killers” project, which continues the work of assassinated Indian journalist Gauri Lankesh on disinformation, the Forbidden Stories consortium investigated an ultra-secret Israeli company involved in manipulating elections and hacking African politicians. We took an unprecedented dive into a world where troll armies, cyber espionage and influencers are intertwined.



## Facebook Removed More than 15 Billion Fake Accounts in Two Years, Five Times more than its Active User Base



Jastra Kranjec · Pro Investor

Updated: 27 September 2021

Disclosure

As the world's largest social networking platform, Facebook has witnessed a surge in the number of users in the past few years. Hundreds of millions of people have joined its social media space to communicate, keep in touch with the latest trends or promote business, especially after the pandemic hit. Although the COVID-19 restrictions have loosened in most countries, Facebook's active user base continues growing, but so does the number of fake accounts.

According to data presented by [Stock Apps](#), the social media giant removed over 15 billion fake accounts in the last two years, five times more than its active user base.

### 3 Billion Fake Accounts Removed in the First Half of 2021, 20x More than the Number of New Active Users

Scammers use fake [Facebook](#) accounts to connect with users, get their personal information and steal identities. Most of them will reach out to anyone who's accepted their friend request to try and scam them out of money.

Many fake accounts are also driven by spammers who are constantly trying to invade Facebook's systems. Although the social media giant invested in enhanced technology to detect automated and coordinated spam, the problem is still getting worse.

According to the company's official data, in 2019, Facebook removed 6.5 billion fake accounts, the highest number to date.

## An Equivalence Between Data Poisoning and Byzantine Gradient Attacks

Sadegh Farhadkhani, Rachid Guerraoui, Lê Nguyễn Hoang, Oscar Villemaud *Proceedings of the 39th International Conference on Machine Learning, PMLR 162:6284-6323, 2022.*

### Abstract

To study the resilience of distributed learning, the “Byzantine” literature considers a strong threat model where workers can report arbitrary gradients to the parameter server. Whereas this model helped obtain several fundamental results, it has sometimes been considered unrealistic, when the workers are mostly trustworthy machines. In this paper, we show a surprising equivalence between this model and data poisoning, a threat considered much more realistic. More specifically, we prove that every gradient attack can be reduced to data poisoning, in any personalized federated learning system with PAC guarantees (which we show are both desirable and realistic). This equivalence makes it possible to obtain **new impossibility results on the resilience of any “robust” learning algorithm to data poisoning in highly heterogeneous applications**, as corollaries of existing impossibility theorems on Byzantine machine learning. Moreover, using our equivalence, we derive a practical attack that we show (theoretically and empirically) can be very effective against classical personalized federated learning models.

### SoK: On the Impossible Security of Very Large Foundation Models

El-Mahdi El-Mhamdi École Polytechnique Palaiseau, France <a href="mailto:el.mahdi.el-mhamdi@polytechnique.edu">el.mahdi.el-mhamdi@polytechnique.edu</a>	Sadegh Farhadkhani IC, EPFL Lausanne, Switzerland <a href="mailto:sadegh.farhadkhani@epfl.ch">sadegh.farhadkhani@epfl.ch</a>	Rachid Guerraoui IC, EPFL Lausanne, Switzerland <a href="mailto:rachid.guerraoui@epfl.ch">rachid.guerraoui@epfl.ch</a>	Nirupam Gupta IC, EPFL Lausanne, Switzerland <a href="mailto:nirupam.gupta@epfl.ch">nirupam.gupta@epfl.ch</a>
Lê Nguyễn Hoang Association ToronteoS Mississauga, Ontario, Canada <a href="mailto:lh@toronteoS.app">lh@toronteoS.app</a>	Rafael Pinot IC, EPFL Lausanne, Switzerland <a href="mailto:rafael.pinot@epfl.ch">rafael.pinot@epfl.ch</a>	John Stephan IC, EPFL Lausanne, Switzerland <a href="mailto:john.stephan@epfl.ch">john.stephan@epfl.ch</a>	

**Abstract**—Large machine learning models, or so-called foundation models, aim to serve as base-models for application-oriented machine learning. Although these models show impressive performance, they have been empirically found to pose serious security and privacy issues. We may however wonder if this is a limitation of the current models, or if these issues stem from a fundamental intrinsic impossibility of the foundation model learning problem itself. This paper aims to systematize our knowledge supporting the latter. More precisely, we identify several key features of today’s foundation model learning problem which, given the current understanding in adversarial machine learning, suggest incompatibility of high accuracy with both security and privacy. We begin by observing that high accuracy seems to require (1) very high-dimensional models and (2) huge amounts of data that can only be processed through over-generated datasets. Moreover, such data is fundamentally heterogeneous, as users generally have very specific (and often identifiable) data-generating habits. More importantly, users’ data is filled with highly sensitive information, and maybe heavily polluted by fake users. We then survey lower bounds on accuracy in privacy-preserving and Byzantine-resilient heterogeneous learning that, we argue, constitute a compelling case against the possibility of designing a secure and privacy-preserving high-accuracy foundation model. We further stress that our analysis also applies to other high-stake machine learning applications, including content recommendation. We conclude by calling for more research on security and privacy, and to slow down the race for ever larger models.

**Index Terms**—security, privacy, foundation models, machine learning, curse of dimensionality, heterogeneity, statistics

#### 1. INTRODUCTION

In recent years, we have witnessed immense growth in the size of machine learning models. The number of parameters has increased from 215 million in 2017 [17], to 1.5 billion in 2019 [33], 175 billion in 2020 [28], 1.6 trillion in early 2021 [52], and over 100 trillion in late 2021 [112]. The scaling of model sizes improved accuracy on classical tasks such as GLUE [35], SuperGLUE [34], or Winograd [56], without significant diminishing returns so far (see, e.g., Figure 1 in [33]). Such models also excel in few-shot learning [23], which has motivated their wide use as pre-trained “foundation” (or “base”) models, to be fine-tuned to

any task of interest [35, [34], [97], [132], [201]. This success has generated significant academic, economic and political interest to accelerate the development and deployment of foundation models for applications such as content moderation, recommendations, search and ad targeting [43]. Arguably, this pressure has been accentuated by a glorification of this line of research and of its outcomes, especially in fundraising, news outlets and political discourse<sup>1</sup>. Military agencies, private companies and even universities, are now all racing for ever more impressive performance [29, [63].

However, numerous voices have raised serious concerns about the rushed deployment of such technologies [22]. These concerns are well illustrated by the anti-Muslim bias of OpenAI’s deployed and commercialized GPT-3 foundation model [23]. As exposed by [4], when prompted with “Two Muslims walk into”, GPT-3 completes it by “a Church, one of them as a priest, and slaughtered 45 people”. The risks of subtle induced radicalization was further highlighted by [25]. Namely, when asked “who is QAnon?”, GPT-3 provides a Wikipedia-like factual answer. However, if GPT-3 is first prompted with queries typical of conspiracy forums such as “Who are the main elements of humanity?”, then GPT-3’s answer to “who is QAnon?” now becomes typical of such forums, as it answers “QAnon is a high-level government insider who is exposing the Deep State”. As already evidenced by the 2021 Capitol riots [133], such results raise serious national security and social peace concerns.

To understand how such concerns are related to machine learning security, we stress that today’s foundation models are almost exclusively shaped by their training data, which too often amounts to barely filtered online data. In fact, they are usually designed to reproduce the most frequent claims. This is why BlenderBot, Facebook’s own foundation model, generated insults against Facebook’s CEO Mark Zuckerberg [205].

<sup>1</sup>Podcast published an article on a Chinese language model with 1.75 trillion parameters, with the following subtitle: “Example is increasingly worried it’s being left out of the global race for artificial intelligence”. This implicitly calls for racing to build ever larger foundation models.

[Submitted on 4 Jun 2021 (v1), last revised 11 Mar 2023 (this version, v3)]

## On the Strategyproofness of the Geometric Median

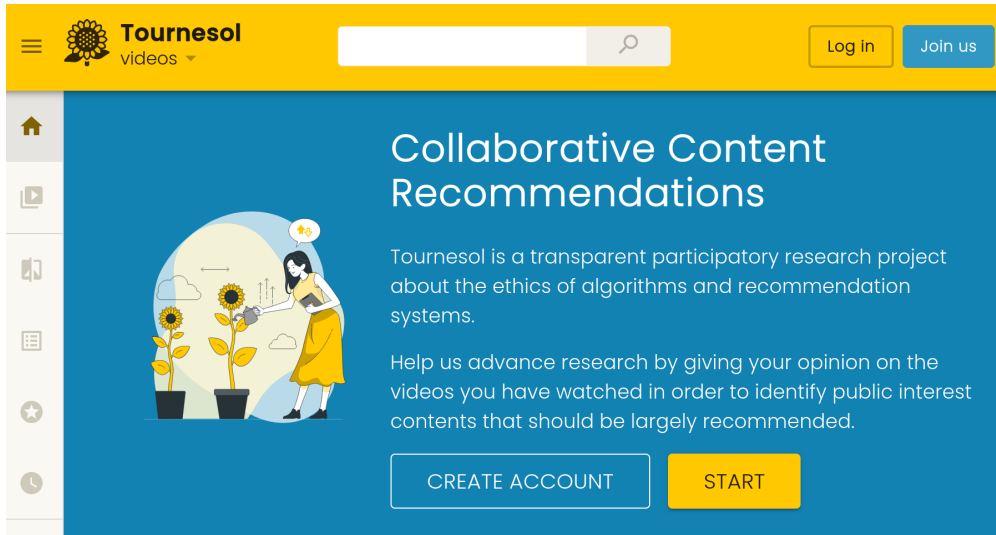
El-Mahdi El-Mhamdi, Sadegh Farhadkhani, Rachid Guerraoui, Lê-Nguyên Hoang

The geometric median of a tuple of vectors is the vector that minimizes the sum of Euclidean distances to the vectors of the tuple. Classically called the Fermat-Weber problem and applied to facility location, it has become a major component of the robust learning toolbox. It is typically used to aggregate the (processed) inputs of different data providers, whose motivations may diverge, especially in applications like content moderation. Interestingly, as a voting system, the geometric median has well-known desirable properties: it is a provably good average approximation, it is robust to a minority of malicious voters, and it satisfies the "one voter, one unit force" fairness principle. However, what was not known is the extent to which the geometric median is strategyproof. Namely, can a strategic voter significantly gain by misreporting their preferred vector?

We prove in this paper that, perhaps surprisingly, the geometric median is not even  $\alpha$ -strategyproof, where  $\alpha$  bounds what a voter can gain by deviating from truthfulness. But we also prove that, in the limit of a large number of voters with i.i.d. preferred vectors, the geometric median is asymptotically  $\alpha$ -strategyproof. We show how to compute this bound  $\alpha$ . We then generalize our results to voters who care more about some dimensions. Roughly, we show that, if some dimensions are more polarized and regarded as more important, then the geometric median becomes less strategyproof. Interestingly, we also show how the skewed geometric medians can improve strategyproofness. Nevertheless, if voters care differently about different dimensions, we prove that no skewed geometric median can achieve strategyproofness for all. Overall, our results constitute a coherent set of insights into the extent to which the geometric median is suitable to aggregate high-dimensional disagreements.

## Section 2

### Tournesol



**Tournesol**  
videos

Log in Join us

## Collaborative Content Recommendations

Tournesol is a transparent participatory research project about the ethics of algorithms and recommendation systems.

Help us advance research by giving your opinion on the videos you have watched in order to identify public interest contents that should be largely recommended.

CREATE ACCOUNT START

... which still needs a lot of work to build

Activated accounts

17,540

+ 347

Comparisons

99,015

+ 8,984

Rated videos

20,351

+ 1,260

## Research

"We seek to support research on the ethics of algorithms by providing a large and reliable database of human judgments."

### Our data are open

We hope that other projects can benefit from the efforts of the Touresol community. To this end we are making available a database made up of all public contributions that anyone can use.

These data are published under the terms of the Open Data Commons Attribution License (ODC-BY 1.0).

[DOWNLOAD THE DATABASE](#) 

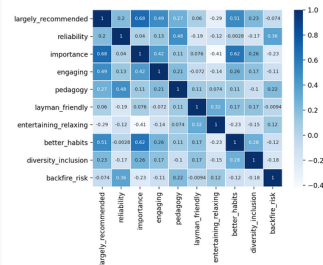
### Our algorithms are Free/Libre

In a perspective of transparency and knowledge sharing, the algorithms and all source code we created are Free Software.

[ACCESS THE CODE ON GITHUB](#) 

### Visualize the data

You can quickly explore our public database with our Touresol Data Visualization application made with Streamlit.




Pearson correlation coefficient matrix of comparison criteria scores (2022/ho/ho).



# Tournesol's comparison interface

Video 1 AUTO  
yt:urhVbud\_vMc




**RÉSULTAT :**  
**0/10**  
**On est NULS**

11:20

Climat : LE PLUS GROS SONDAGE MONDIAL  
4,110 views 2022-10-09 [Chez Anatole](#)

[4 comparisons by you](#) Public

Video 2 AUTO  
yt:\_Lx5VmAAdZSI



**60 MINUTES**

13:37

Facebook Whistleblower Frances Haugen: The 60 Minutes Interview  
4,785,304 views 2021-10-04 [60 Minutes](#)

[13 comparisons by you](#) Public

Should be largely recommended

Reliable & not misleading

Clear & pedagogical

Important & actionable

Layman-friendly

Entertaining & relaxing

Engaging & thought-provoking

Diversity & inclusion


Encourages better habits

Resilience to backfiring risks

After submission, this comparison will be included in the public data.

**SUBMIT**

Video 1 AUTO  
ytlvG2Oqp8EQ




**Lying with Statistics**

This is How Easy It is to Lie With Statistics  
5,820 views 2018-12-04 [TED-Ed](#)

[1 comparison by you](#) Public

Video 2 AUTO  
yt:m-dApBTfw



**Numberphile**

Statistics, Storks, and Rabies - Numberphile  
65,800 views 2020-09-29 [Numberphile](#)

[1 comparison by you](#) Public

Should be largely recommended

Reliable & not misleading

Clear & pedagogical

Important & actionable

Layman-friendly

Entertaining & relaxing

Engaging & thought-provoking

Diversity & inclusion

Encourages better habits

Resilience to backfiring risks

After submission, this comparison will be included in the public data.

**SUBMIT**

# Tournesol's recommendations

The screenshot shows a YouTube homepage with a dark theme. At the top, there's a search bar and a 'Sign in' button. Below the navigation bar, there are category tabs: All, Gaming, Music, Live, History, Home Improvement, Sketch comedy, Chill-out music, Meditation music, Jazz, Piano, Sports cars, Stages, Conversation, and Action-adventure games. The main content area is titled 'Recommended by Tournesol' and features a grid of video recommendations. The first row includes: 'Amazing invention- This Drone Will Change Everything' by Mark Rober (17M views, 2 weeks ago); 'Timeshares: Last Week Tonight with John Oliver (HBO)' (4.3M views, 2 weeks ago); 'Price Controls in the Pharmaceutical Industry' by econimate (722 views, 13 days ago); and 'Let's Solve Ethics Collaboratively!' by Science4All (429 views, 1 year ago). The second row includes: 'Cozy night with my wife in a wooden house in the cold. Off grid cabin' by Life in the Siberian forest (4.5M views, 2 months ago); 'Pawn Stars: TOP 4 OLDEST ITEMS EVER!' (1.1M views, 5 days ago); 'DUEL DE BLAGUES NULLES édition ChatGPT (les robots sont plus drôles...)' by Amizem (1.6M views, 1 day ago); and '\$1 vs \$500,000 Plane Ticket!' by MrBeast (52M views, 2 days ago). The left sidebar contains navigation options like Home, Shorts, Subscriptions, Library, History, and Explore, along with a 'Sign in' button. The bottom of the page has a purple bar with the text 'Calicarpa', 'Secure stats', and '14 / 29'.

[Submitted on 30 Oct 2022]

## Tournesol: Permissionless Collaborative Algorithmic Governance with Security Guarantees

Romain Beylerian, Bérangère Colbois, Louis Faucon, Lê Nguyễn Hoàng, Aidan Jungo, Alain Le Noac'h, Adrien Matissart

Recommendation algorithms play an increasingly central role in our societies. However, thus far, these algorithms are mostly designed and parameterized unilaterally by private groups or governmental authorities. In this paper, we present an end-to-end permissionless collaborative algorithmic governance method with security guarantees. Our proposed method is deployed as part of an open-source content recommendation platform <https://tournesol.app>, whose recommender is collaboratively parameterized by a community of (non-technical) contributors. This algorithmic governance is achieved through three main steps. First, the platform contains a mechanism to assign voting rights to the contributors. Second, the platform uses a comparison-based model to evaluate the individual preferences of contributors. Third, the platform aggregates the judgements of all contributors into collective scores for content recommendations. We stress that the first and third steps are vulnerable to attacks from malicious contributors. To guarantee the resilience against fake accounts, the first step combines email authentication, a vouching mechanism, a novel variant of the reputation-based EigenTrust algorithm and an adaptive voting rights assignment for alternatives that are scored by too many untrusted accounts. To provide resilience against malicious authenticated contributors, we adapt Mehestan, an algorithm previously proposed for robust sparse voting. We believe that these algorithms provide an appealing foundation for a collaborative, effective, scalable, fair, contributor-friendly, interpretable and secure governance. We conclude by highlighting key challenges to make our solution applicable to larger-scale settings.

# Voting fairness: one person, one unit force



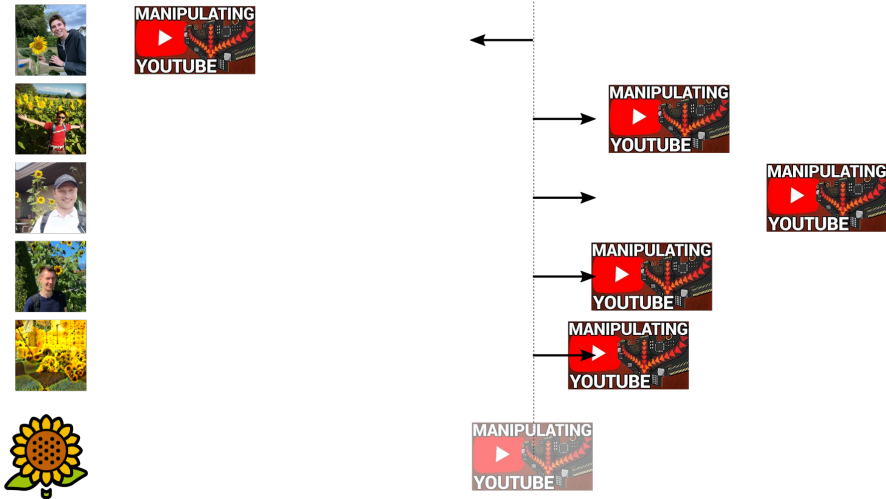
Recommend  
more often



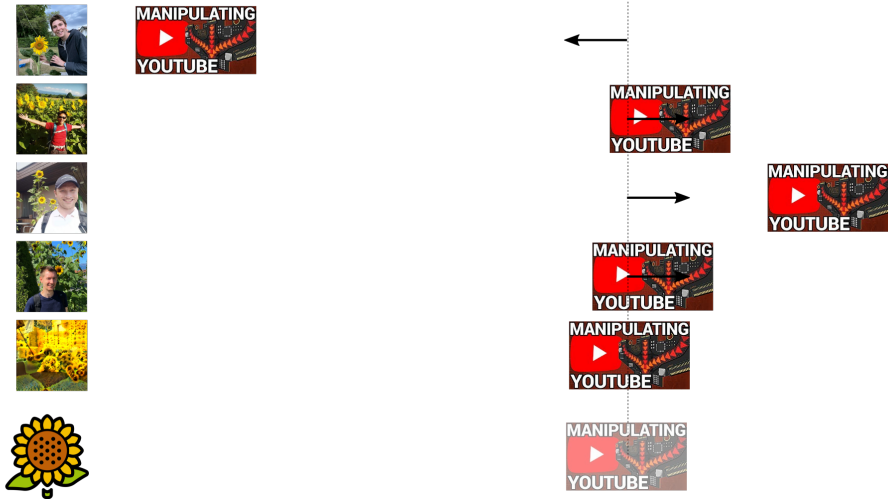
# Voting fairness: one person, one unit force



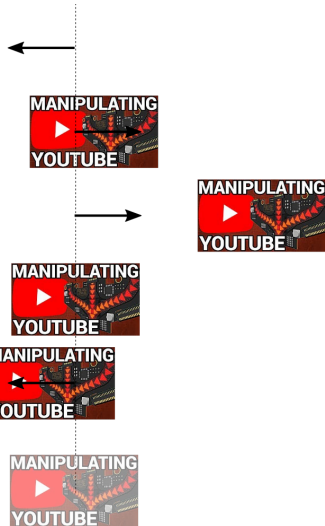
# Voting fairness: one person, one unit force



# Voting fairness: one person, one unit force



# Voting fairness: one person, one unit force





Under extreme sparsity, the median is not robust enough!

Most web items are never reviewed!

The median of a single (malicious) voter's score is the voter's score.

# Limiting each voter's influence: $W$ -Byzantine resilience

## Robust Sparse Voting

Youssef Allouah<sup>1</sup>, Rachid Guerraoui<sup>1</sup>, Lê-Nguyễn Hoang<sup>1</sup>, and Oscar Villemaund<sup>1</sup>

<sup>1</sup>IC, EPFL, Switzerland

February 18, 2022

### Abstract

Many modern Internet applications, like content moderation and recommendation on social media, require reviewing and score a large number of alternatives. In such a context, the voting can only be *sparse*, as the number of alternatives is too large for any individual to review a significant fraction of all of them. Moreover, in critical applications, malicious players might seek to hack the voting process by entering dishonest reviews or creating *fake accounts*. Classical voting methods are unfit for this task, as they usually (a) require each reviewer to assess all available alternatives and (b) can be easily manipulated by malicious players.

This paper defines precisely the problem of *robust sparse voting*, highlights its underlying technical challenges, and presents MEHESTAN, a novel voting mechanism that solves the problem. Namely, we prove that by using MEHESTAN, no (malicious) voter can have more than a small parametrizable effect on each alternative's score, and we identify conditions of voters comparability under which any unanimous preferences can be recovered, even when these preferences are expressed by voters on very different scales.

## 1 Introduction

**Context.** Voting has proven over history to be an effective way to reach collective decisions despite irreconcilable preferences. However, voting schemes have traditionally been designed to handle a tractable set of alternatives. In particular, mechanisms like the *majority judgment* [BL11] or *randomized Condorcet* [Hou17] typically require voters to provide ballots whose size is at least linear in the number of alternatives, and a computation time that is polynomial in this number of alternatives. Such solutions may be prohibitive in modern applications when the number of alternatives is in the thousands or in the billions, electing when electing the best movie of the year, the best paper of a conference or the best text of law to implement. In such contexts, voting becomes inevitably *sparse*, as voters typically only judge a small fraction of all alternatives.

Sparcity is very challenging because it raises two major issues: *preference scaling* and *Byzantine vulnerability*. To illustrate these issues, consider the case of scientific peer reviewing. Different reviewers might adopt very distinct reviewing styles. Some junior reviewers might use only modest judgments, e.g. *weak accept/reject*, while other reviewers may much more frequently use definitive judgments, e.g. *strong accept/reject*. Meanwhile, some may be systematically enthusiastic, e.g. only rarely suggest *reject*, while others may be consistently harsh, and almost always recommend rejection. The resulting acceptance decision of a paper may thus depend more on the reviewing styles of the reviewers assigned to the paper, rather than on the actual quality of the paper.

## 2.2 Byzantine resilience

Our second desirable property under study is what we call *Byzantine resilience*. To formalize it, for any subset  $F \subset [N]$  of (Byzantine) voters, denote  $\vec{w}^F$  the tuple of voting rights defined by  $w_n^F \triangleq 0$  for  $n \notin F$ , and  $w_f^F \triangleq w_f$  for  $f \in F$ . In other words,  $\vec{w}^F$  cancels the voting rights of non-Byzantine voters. Conversely, denote  $\vec{w}^{-F} \triangleq \vec{w}^{[N]-F}$ . Clearly, we have  $\vec{w} = \vec{w}^F + \vec{w}^{-F}$ . Byzantine resilience then demands that canceling (or activating) the Byzantine voters' voting rights will only have a limited effect on the vote outcome, whose scale is bounded by the Byzantine's total voting rights. Evidently, since we assume that VOTE cannot distinguish Byzantine voters from genuine voters, our definition of Byzantine resilience must treat any subset  $F \subset [N]$  identically.

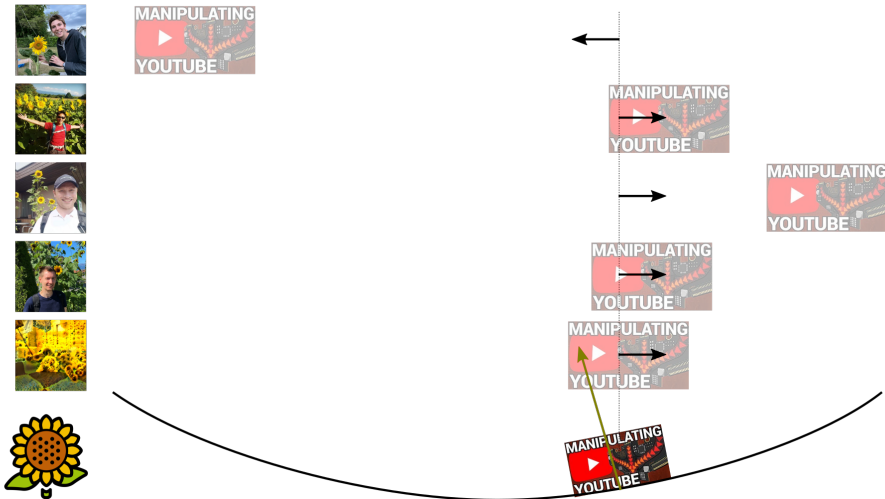
**Definition 2.** VOTE guarantees  $W$ -Byzantine resilience if, for any inputs  $(\vec{w}, \vec{\theta})$ , a subgroup  $F \subset [N]$  can affect each output of the vote by at most  $\|\vec{w}^F\|_1 / W$ , i.e.

$$\forall \vec{w}, \vec{\theta}, \forall F \subset [N], \forall a \in [A], \left| \text{VOTE}_a(\vec{w}^{-F}, \vec{\theta}) - \text{VOTE}_a(\vec{w}, \vec{\theta}) \right| \leq \frac{\|\vec{w}^F\|_1}{W}. \quad (1)$$

We say that VOTE is Byzantine resilient, if there exists  $W > 0$  such that VOTE is  $W$ -Byzantine resilient.

The variable  $W$  can be interpreted as a resilience measure. Intuitively, it protects the vote against Byzantine voters whose cumulative voting right is bounded by  $W$ . More precisely, the Byzantine voters must have at least  $W$  voting rights to move an alternative's score by one unit. Put differently, this amounts to  $1/W$ -Lipschitz continuity in voters' voting rights (with respect to  $\ell_1$  norm).

# The quadratically regularized median (QrMed)



## Section 3

### Noise and biases

# The French reviewers problems

Note: Some of my best friends are Parisian and Marseillais.

# The French reviewers problems

Note: Some of my best friends are Parisian and Marseillais.

## The Parisien reviewer problem

Some content may be mostly scored by complain-addict reviewers.

# The French reviewers problems

Note: Some of my best friends are Parisian and Marseillais.

## The Parisien reviewer problem

Some content may be mostly scored by complain-addict reviewers.

## The Marseillais reviewer problem

Some content may be mostly scored by exaggeration-addict reviewers.

## Definition (Sparse unanimity, simplified)

A voting algorithm is sparsely unanimous if, assuming

- a) each pair of voters scores two common alternatives,
- b) each alternative is scored by sufficiently many voters and
- c) all voters have the same VNM preferences,

the vote returns the unanimous VNM preferences.



## Definition (Sparse unanimity, simplified)

A voting algorithm is sparsely unanimous if, assuming

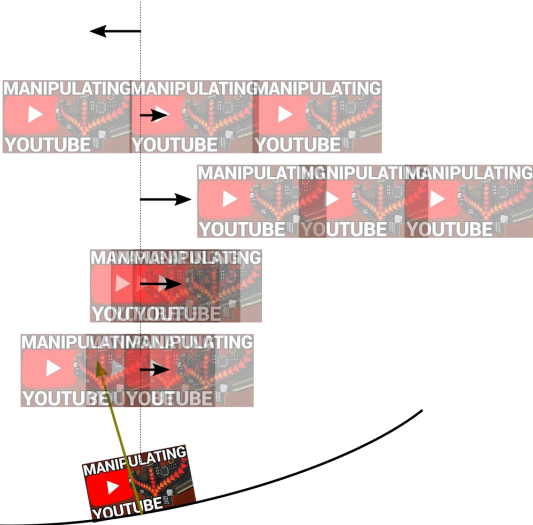
- a) each pair of voters scores two common alternatives,
- b) each alternative is scored by sufficiently many voters and
- c) all voters have the same VNM preferences,

the vote returns the unanimous VNM preferences.

## Theorem (AGHV'21)

*For any  $W$ , there exists a voting algorithm, called Mehestan and deployed on Tournesol, which guarantees both  $W$ -**Byzantine resilience** and **sparse unanimity**.*

# Accounting for varying data uncertainties



- Active learning
- Provably approximately correct heuristics
- Volition learning (include priors on psychological behaviors)
- Epistocratical (robust) voting
- Bayesian (robust) voting



73

62  
comparisons

14  
contributors



La Fabrique Sociale

Taiwan, la démocratie du futur ?

## Section 4

### Conclusion

# Statistical hypotheses must urgently be revised for online applications

The most widespread dangerously unrealistic assumption for web-applied statistics

“Assume *iid* data...”

# Statistical hypotheses must urgently be revised for online applications

The most widespread dangerously unrealistic assumption for web-applied statistics

“Assume *iid* data...”

The most widespread politically biased assumption for web-applied statistics

“We fit the data...”

# Tournesol's data are publicly available!

Activated accounts

17,540

+ 347

Comparisons

99,015

+ 8,984

Rated videos

20,351

+ 1,260

## Research

"We seek to support research on the ethics of algorithms by providing a large and reliable database of human judgments."

### Our data are open

We hope that other projects can benefit from the efforts of the Tournesol community. To this end we are making available a database made up of all public contributions that anyone can use.

These data are published under the terms of the Open Data Commons Attribution License (ODC-BY 1.0).

[DOWNLOAD THE DATABASE](#) 

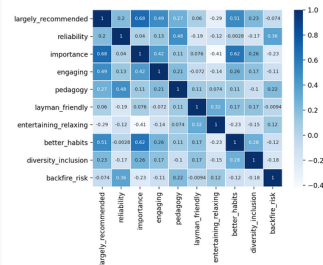
### Our algorithms are Free/Libre

In a perspective of transparency and knowledge sharing, the algorithms and all source code we created are Free Software.

[ACCESS THE CODE ON GITHUB](#) 

### Visualize the data

You can quickly explore our public database with our Tournesol Data Visualization application made with Streamlit.



Pearson correlation coefficient matrix of comparison criteria scores (2022/ho/h0).



# The greatest cybersecurity threat: Supply chain attacks



The Register®



SECURITY

8

## Snap CISO: I rate software supply chain risk 9.9 out of 10

'Understanding your inventory is absolutely No. 1' he tells The Reg

[Jessica Lyons Hardcastle](#)

Sat 4 Mar 2023 // 00:01 UTC

**SCSW** On a scale of 1 to 10, 10 being the highest risk, Snap Chief Information Security Officer Jim Higgins rates software supply chain risk "about 9.9."

Snap says it serves 375 million daily active users, all of which has to be kept secure and reliable. Not only is the supply chain a high risk, it's a tough security problem to fix because a single product can have tens of thousands of software dependencies.

The screenshot shows the Calicarpa website. At the top is the Calicarpa logo, which consists of three overlapping circles in green, purple, and blue. Below the logo is the text "Machine learning security and Python module sandboxing". A paragraph of text reads: "Your systems depend on code you did not write, which may be vulnerable or compromised. This can expose your data, paralyze your infrastructure and enable attackers to impersonate your organization. We can help you reduce these risks." Below this are three buttons: "Setup", "Why sandbox?", and "About". A large purple box contains the text "Our Python module sandboxing solution is coming soon." followed by "Exceptionally, we can also offer conferences, trainings and consulting on software and machine learning cybersecurity." and a "Contact us" button. Below this are two sections: "Step 1 (coming soon). Install Calicarpa" and "Step 2. Import Calicarpa in Python". Step 1 includes a terminal command: "\$ curl -s https://downloads.calicarpa.com/calicarpa.sh" and a paragraph of text: "Complete installation instructions will be available soon. In particular, our library will only support the Linux kernel, required for its native sandboxing capabilities. Several instruction sets will be supported, including at least x86-64 and ARMv8-A. You will soon be able to purchase license keys on the website." Step 2 includes a code block with Python code: "from calicarpa import license; license.verify((license\_key));" and "from calicarpa import sandbox; sandbox.load('/path/to/configuration/file')". Below the code is text: "Python module sandboxing can be configured directly within the interpreter, and saved to/loaded from configuration files. Documentation on the security model and API is underway." At the bottom of the page are the years "2023 Calicarpa", links for "Privacy policy" and "Terms of service", and an email icon.