# Position: Generative AI Regulation Can Learn from Social Media Regulation

Ruth E. Appel<sup>1</sup>

## Abstract

There is strong agreement that generative AI should be regulated, but strong disagreement on how to approach regulation. While some argue that AI regulation should mostly rely on extensions of existing laws, others argue that entirely new laws and regulations are needed to ensure that generative AI benefits society. In this position paper, we argue that the debates on generative AI regulation can be informed by evidence on social media regulation. For example, AI companies have faced allegations of political bias which resemble the allegations social media companies have faced. First, we compare and contrast the affordances of generative AI and social media to highlight their similarities and differences. Then, we discuss four specific policy recommendations based on the evolution of social media and their regulation: (1) counter bias and perceptions thereof (e.g., via transparency, oversight boards, researcher access, democratic input), (2) address specific regulatory concerns (e.g., youth wellbeing, election integrity) and invest in trust and safety, (3) promote computational social science research, and (4) take on a more global perspective. Applying lessons learnt from social media regulation to generative AI regulation can save effort and time, and prevent avoidable mistakes.

## 1. Introduction

When Google's generative AI model Gemini produced images of racially diverse Nazis in early 2024, it led to a public outcry and allegations of anti-conservative bias (Robertson, 2024). Almost a decade earlier, the first allegations of anti-conservative bias were made against social media platforms like Facebook (Barrett & Sims, 2021), and they have persisted e.g. during Senate hearings (Romm, 2019) and when President Trump was banned from Twitter (now X) and Facebook (Barrett & Sims, 2021). This shows that the content moderation challenges that emerging technologies face are not entirely new. Media scholars have called attention to the fact that new technologies often elicit similar questions and concerns as their predecessors (Wartella & Reeves, 1985). Generative AI is the latest technology to garner widespread attention and raise societal and regulatory concerns, but so have social media and other technologies before it.

In this paper, we argue that generative AI regulation can learn from social media regulation, which has evolved over the past two decades. While there is strong agreement that generative AI should be regulated—evidenced by the large number of recent regulatory efforts across countries and stakeholders (Zaidan & Ibrahim, 2024)-, there is strong disagreement on how to approach regulation. Some argue that AI regulation should mostly rely on extensions of existing laws (Huttenlocher et al., 2023), while others argue that entirely new laws and regulations are needed and have proposed laws and regulations such as the EU AI Act (European Parliament and Council, 2024), White House Executive Orders on AI (Executive Office of the President, 2023; 2025), or California's vetoed AI Safety Bill SB 1047 (Wiener et al., 2024). Analyzing the evolution of social media regulation can provide insights into which approaches to regulation are promising when it comes to generative AI, which in turn can prevent avoidable mistakes, and save effort and time.

Learning from social media regulation is imperative because AI regulation is urgent. Misused or misaligned generative AI can cause severe harms (Weidinger et al., 2021; Marchal et al., 2024; MITRE Corporation, 2021), and the risks become even greater as generative AI advances (Hendrycks et al., 2025). Yet, despite a patchwork of emerging global regulations, effective global regulation of generative AI is lacking (Alanoca et al., 2025).

Specific learnings for generative AI regulation based on the evolution of social media regulation include investments in (1) efforts to counter bias and perceptions thereof (e.g., via transparency, researcher access, oversight boards, democratic input), (2) specific areas of regulatory concern and trust and safety, (3) computational social science research,

<sup>&</sup>lt;sup>1</sup>Stanford University, Stanford, CA, United States. Correspondence to: Ruth E. Appel <rappel@cs.stanford.edu>.

Proceedings of the  $42^{nd}$  International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

and (4) a more global perspective (see Section 3).

The focus of this paper is on content moderation, i.e. how to design and regulate the content generated by generative AI models or shown on social media platforms, because social media regulation holds particularly relevant insights in this area. Further, the paper focuses on regulation in a broad sense, which can include self-regulation of industry players and formal laws such as the EU AI Act (European Parliament and Council, 2024) or White House Executive Orders on AI (Executive Office of the President, 2023; 2025).

First, we compare and contrast the affordances of generative AI and social media to highlight their similarities and differences. Second, we provide policy recommendations based on the evolution of social media and their regulation. Third, we discuss implementation challenges such as political polarization and the pace of technological change. Fourth, we engage with alternative views before we conclude.

# 2. Affordances of Generative AI and Social Media

To shed light on the similarities and differences between specific media, we can analyze their affordances. For the purposes of this paper, we define affordances as the features that characterize a medium in its relationship to its users (for a detailed discussion of different definitions and the evolution of the term affordances, see McGrenere & Ho, 2000; Ronzhyn et al., 2023). Both generative AI, e.g. in the form of a chatbot like OpenAI's ChatGPT or Anthropic's Claude, and social media, e.g. in the form of Meta's Facebook or X (formerly Twitter), can be considered media that allow to create and distribute content and are shaped by specific features. The features discussed here pertain to a medium in general, but may not apply to every instance, that is, a specific generative AI application or social media platform may differ from the norm in terms of its affordances.

Based on an analysis of commonly used generative AI applications (e.g., ChatGPT and Claude) as well as social media applications (e.g., Facebook and X), we identified key features that generative AI and social media share or that differentiate them. The analysis of features is grounded in work by Clark (1996), who discusses several features of media that fall into three categories: medium, control, and immediacy. Since Clark (1996)'s features focus on the affordances of face-to-face communication,<sup>1</sup> we added new features and removed features that are less relevant to the

comparison of generative AI and social media. We also added the feature of interactivity discussed by Rafaeli & Sudweeks (1997). We will point out each feature that is adapted from Clark (1996) or Rafaeli & Sudweeks (1997).

We will first address why social media is comparable to generative AI in key aspects that have implications for technology regulation. Then, we will engage with differences in the affordances of social media and generative AI to show that the analogy is useful, but imperfect.

#### 2.1. Generative AI and Social Media Are Comparable in Key Aspects

The analogy between generative AI and social media is valuable because both media share key features. Importantly, the shared affordances of generative AI and social media imply that both of these media necessarily moderate content and thus face complex content moderation challenges and public scrutiny.

Table 1 shows key similarities between generative AI and social media when it comes to the features of each medium. Both generative AI and social media allow for spatial separation, that is, the conversation partners usually generate content in different physical spaces-e.g., in a home office and at a data center for generative AI-and are not copresent (copresence is one of the features described in Clark (1996)). Both media feature interactivity and respond interactively to user input, which makes them engaging (Rafaeli & Sudweeks, 1997) (interactivity is defined and discussed in Rafaeli & Sudweeks (1997)). Both generative AI and social media are *recording* user data (the recording feature is adapted from Clark (1996)'s recordlessness feature). Both media can learn about a user's context and their preferences over time for output *personalization*, e.g. by updating the chatbot's memory or personalizing a recommendation algorithm. Further, both generative AI and social media can feature general content, i.e. content on all kinds of domains (e.g., hobbies, jobs, politics). Both are powered by artificial intelligence (AI), that is, they rely on learning patterns from data to perform well on tasks such as generating text or recommending content, although generative AI relies on more recent deep learning models while social media tends to rely on traditional machine learning approaches such as recommender systems. Both media also feature abstraction, that is, they hide the complex technical implementation details from the user behind a simple user interface. Further, generative AI and social media algorithms tend to be blackbox, that is, algorithmic decisions are intransparent-almost always for users, but often also for experts because mechanistic interpretability (Bereska & Gavves, 2024) that can explain why a deep learning model made a certain decision is in its infancy.

With regards to control features (Clark, 1996), both gen-

<sup>&</sup>lt;sup>1</sup>There are contextual differences between face-to-face communication on the one hand and generative AI and social media on the other, such as where and why they may be used. This paper focuses on the comparison of generative AI and social media, and therefore focuses on features in Clark (1996)'s model that are pertinent to generative AI and social media, but not the comparison to other media.

Feature	Definition	Generative AI	Social Media
Medium			
Spatial separation	Content is generated in different locations	Yes	Yes
Direct connection	Medium is conversation partner	Yes	No
User connections	Medium connects user to other users	No	Yes
Interactivity	Medium responds interactively to user input	Yes	Yes
Dialogue-by-default	Actions occur in a dialogue	Yes	No
Recording	User actions are recorded	Yes	Yes
Personalization	User context and preferences are learnt over time	Yes	Yes
Single output	Medium presents usually just a single output	Yes	No
Infinite content	Content is served infinitely	No	Yes
General content	Content can pertain to any domain	Yes	Yes
General purpose	Medium serves many functions	Yes	No
Use of AI	Medium learns patterns from data	Yes	Yes
Abstraction	Medium hides its complexity	Yes	Yes
Black-box	How algorithmic decisions are made is intransparent	Yes	Yes
Control			
Content moderation	Content is moderated at all	Yes	Yes
Invisible content moderation	Most content moderation is not visible to the user	Yes	No
Content moderation pre-generation	Content is moderated before it is received by the user	Yes	No
Self-determination	User can decide themselves how to act	Yes	Yes
Self-expression	User can express themselves	Yes	Yes
Simultaneity	User can receive and produce content concurrently	No	Yes
Immediacy			
Instantaneity	Actions are perceived almost immediately	Yes	Yes
Evanescence	Medium quickly recedes to the background	Yes	Yes

Table 1. Comparison	of affordances o	f generative A	and social media
---------------------	------------------	----------------	------------------

*Note*: The features spatial separation, recording, self-determination, self-expression, simultaneity, instantaneity, and evanescence, as well as the categories medium, control and immediacy are based on Clark (1996). The feature interactivity is based on Rafaeli & Sudweeks (1997). Instances where features of generative AI are similar to features of social media are highlighted in bold.

erative AI and social media feature *content moderation*, that is, the medium shapes what content is allowed to appear. Both media also meet Clark (1996)'s criteria for *self-determination*, i.e. a user's ability to decide themselves how to act, and *self-expression*, i.e. a user's ability to express themselves on a medium.

With regards to *immediacy* (Clark, 1996), both generative AI and social media share *instantaneity* (Clark, 1996), i.e. that actions are perceived almost immediately, and *evanescence* (Clark, 1996), i.e. that the medium recedes to the background quickly once it is not actively used anymore.

Beyond features, the evolution of generative AI is similar to the evolution of social media in that both are characterized by limited, lagging regulation and large inflows of funding for technology entrepreneurship in this space (Stern, 2023).

#### 2.2. Generative AI and Social Media Are Not Perfectly Comparable

While the shared affordances highlight the value of comparing generative AI to social media, we acknowledge that the analogy is imperfect. By definition, an analogy is not a perfect match. As Jacob Stern put it: "[T]his is just the nature of analogies: They are illuminating but incomplete" (Stern, 2023). Table 1 reveals differences in affordances between generative AI and social media. With regards to features of the medium, generative AI and social media show some variation. While generative AI such as ChatGPT constitutes a conversation partner that is in *direct connection* with the user, social media foster user connections—connections between users. These differences in connection also imply that generative AI tends to be more private by default, since conversations are rarely shared with other users. Whereas generative AI interacts in a dialogue-by-default manner with the user, social media is merely mediating between the user and their human conversation partners (e.g., when a social media algorithm displays one user's post on another user's feed) and tend to involve a sequence of one-off actions. While generative AI tends to respond to prompts, usually with a single output instead of multiple outputs, and does not continue to serve content unless the user requests it, social media often feature infinite content via mechanisms such as infinite scroll (Sharma & Murano, 2020) or autoplay (Lukoff et al., 2021), which serve content as long as the user is on the platform and encourage passive rather than active consumption. The purpose of social media tends to be focused on social communication, while generative AI is considered a general purpose technology that could serve various functions, including as a text writer or reviewer, a

calculator, a programmer and much more.

With regards to control features, a feature Clark (1996) proposed is simultaneity, which is the user's ability to receive and produce content concurrently. Simultaneity is given for social media-e.g., one user might send a message at the same time as another user is sending them a message—, but not for generative AI, which operates in a sequential dialogue of user input and model output. Important differences between generative AI and social media are related to content moderation: Even though both generative AI and social media feature content moderation, content moderation in generative AI tends to use invisible content moderation more than social media. Social media platforms may occasionally take hardly visible actions such as downranking posts, but many social media content moderation actions such as removal of a post or user are clearly visible. Generative AI models, on the other hand, are built and fine-tuned to moderate content in a certain way (e.g., to avoid providing dangerous information), without the user necessarily becoming aware of the moderation. Generative AI content moderation may be invisible to the user because the model will usually respond, and not necessarily provide a reason if it refuses to respond to a prompt directly, which makes moderation less obvious than a missing response or a refused response citing the reason for refusal. Relatedly, generative AI models tend to moderate before the content is shown to the user, e.g. by refusing to reply to a prompt, while social media content moderation tends to occur only after content made it onto a platform, e.g. when a post was reported as harmful misinformation.

Beyond specific features of generative AI and social media, there are differences in their context and potential consequences. In terms of business model, most social media companies rely on revenue from advertisements (Center for Humane Technology, 2021), while prominent generative AI companies have so far leaned towards freemium (Kumar, 2014) subscription models. While the potential harm of social media to democracy and society has been an important focus of scholarly and public attention (Persily & Tucker, 2020), some argue that the destructive potential of AI may be at another level since it may present a larger threat (Bostrom, 2013) or stronger geopolitical advantage (Stern, 2023). Generative AI and social media differ also in the level of uncertainty they bring. For example, auditing and discovering vulnerabilities in systems that are probabilistic (Cattell et al., 2024), like generative AI models, implies new complexities that traditional, deterministic social media algorithms do not entail. Finally, generative AI and social media may differ in areas that have so far remained legally uncertain, such as questions of liability (e.g., for harms resulting from media use) and copyright. This means the learnings for generative AI regulation should be based on, and not go beyond key shared features.

# **3.** Learnings from Social Media Regulation for Generative AI Regulation

As the review of the affordances has shown, generative AI and social media share important features, including the use of AI and content moderation. Although generative AI and social media differ on some dimensions, these differences are mostly differences in degree, and not differences in kind when it comes to regulation. Thus, lessons learnt from social media regulation are relevant to generative AI regulation. This paper provides four policy recommendations for generative AI regulation based on the evolution of social media regulation: (1) counter bias and perceptions thereof (e.g., via transparency, oversight boards, researcher access, democratic input), (2) address specific regulatory concerns (e.g., youth wellbeing, election integrity) and invest in trust and safety, (3) promote computational social science research, and (4) take on a more global perspective. Figure 1 provides an overview of these recommendations.

#### 3.1. Counter Bias and Perceptions Thereof

Given that both generative AI and social media share key features—use of content moderation, use of AI, blackbox nature, abstraction of the complexity of algorithmic decision-making such that much of the decision-making is intransparent—, it is no surprise that both generative AI companies and social media companies have faced allegations of bias, including allegations of anti-conservative political bias (Robertson, 2024; Barrett & Sims, 2021). While there is no evidence of anti-conservative bias for social media (Barrett & Sims, 2021), multiple studies have shown political bias in generative AI. For example, compared to representative opinion polls, large language models were found to output biased opinions (Durmus et al., 2023; Santurkar et al., 2023), and multiple studies showed left-leaning bias in generative AI models (Rozado, 2023; Röttger et al., 2024).

Generative AI models have also been shown to exhibit other forms of bias, such as anti-Muslim bias (Abid et al., 2021), bias towards Western culture (Naous et al., 2023), and stereotypical depictions of race, gender, age, nationality, and socioeconomic status (Nangia et al., 2020). Similarly, generative AI models tend to show social identity biases similar to humans (Hu et al., 2024).

Addressing such biases is as important as it is challenging. It is important to address biases because biases can harm and manipulate users. For example, political bias in generative AI models can influence users' opinions (Bai et al., 2023; Jakesch et al., 2023; Matz et al., 2024; Williams-Ceci et al., 2024; Potter et al., 2024; Anthropic, 2024) and decisions (Fisher et al., 2024). Biases may also lead to lower-quality output, entrench historical biases and stereotypes, and undermine trust. It is challenging to address biases because

Generative AI Regulation Can Learn from Social Media Regulation

Policy Recommendation	Key Strategies	Proposed Generative Al Regulation Measures	Social Media Precedents
Counter bias or perceptions thereof	Transparency, oversight, researcher access, democratic input, personalization	Researcher API access, transparency requirements, decentralization	Meta Transparency Center, TikTok Research API, Mastodon's decentralized content moderation
Address specific regulatory concerns and invest in trust and safety	Investment in promoting youth wellbeing, election integrity, and misinformation prevention	Dedicated trust and safety teams, deceptive campaign monitoring	Trust and safety teams at social media companies such as Google and Meta
Promote computational social science research	Multidisciplinary study of platform impact, evaluation of interventions	Al user experience research, interdisciplinary hiring, rigorous impact evaluation	Facebook and Instagram Election Study, multidisciplinary in-house research teams
Take a more global perspective	Local expertise, international hiring, multilingual content moderation	Ensuring safety and performance in diverse contexts, regionally adapted safety policies	Regional regulatory adaptation, Christchurch Call

Figure 1. Policy recommendation overview. Overview of the lessons generative AI regulation can learn from social media regulation.

they are challenging to measure accurately. For example, bias evaluations may be sensitive to the specific prompt design (Röttger et al., 2024) and order effects (Dominguez-Olmedo et al., 2024). Further, it is not clear where exactly biases stem from. Biases can arise at different points in the development and deployment of generative AI, including training and data curation, fine-tuning, evaluation and feedback, real-time moderation, customization and control of models (Suresh & Guttag, 2021; Ferrara, 2023).

Social media companies have taken different approaches to address biases or perceptions thereof that mainly focus on transparency about algorithms and decision-making, gathering input from users and learning from case studies, and increasing user choice.

# 3.1.1. INCREASE TRANSPARENCY AND RESEARCHER ACCESS

The shared features content moderation, use of AI, blackbox and abstraction give rise to transparency challenges for social media and generative AI. Generative AI transparency is lacking as shown by the Foundation Model Transparency Index (Bommasani et al., 2023a; 2024). Social media companies have pursued multiple different approaches to increase transparency and generative AI can learn from this playbook. For example, Facebook's parent company Meta introduced features such as "Why am I seeing this ad?" that allowed users to understand why they were served certain ad content (Thulasi, 2019), created blog posts and a Transparency Center providing some information on the role of AI and other factors in content recommendation (Clegg, 2023; Meta, 2024a), and established an independent oversight board of experts that adjudicates particularly contentious content moderation decisions (Meta, 2024b). These initiatives do not come without problems. In response to the launch of Facebook's oversight board, "The Real Facebook Oversight Board" was created, which brought experts together to argue for more independence, transparency

and regulation (The Real Facebook Oversight Board, 2022). Company policies are also not guaranteed to be permanent. In January 2025, Meta starkly shifted its content moderation policy, limiting its efforts to reduce misinformation and harmful speech and ending a fact-checking program that had provided some transparency about the content circulating on the platform (Isaac & Schleifer, 2025; Iyer, 2025).

An important aspect of transparency is allowing for thirdparty evaluations. Efforts to create research platforms or APIs accessible to researchers, such as the Meta Researcher Platform (Li et al., 2022) and the TikTok Research API (TikTok, 2025), or to design academic-industry collaboration such as the Facebook and Instagram Election Study (Clegg & Nayak, 2020) are helpful, but imperfect (Wagner, 2023). The Coalition for Independent Technology Research was founded after researchers at different institutions faced difficulty maintaining or gaining access to social media data for research purposes (Coalition for Independent Technology Research, 2022). Importantly, we can learn from these shortcomings. Researcher access programs to evaluate technology should be characterized by sufficient resources (including staffing, infrastructure, and funding), incentives that are compatible with academic research (e.g., data retention policies, persistent API access and publication permission for researchers), sound knowledge sharing processes between internal and external researchers to help understand data availability and analysis feasibility, helpful documentation, privacy preserving measures (e.g., aggregation of user data) and timeliness in terms of data access, publication review and addressing issues that researchers discovered. To protect researchers involved, researchers have called for "safe harbors," that is, legal protection for researchers pursuing legitimate research purposes, initially for social media (Abdo et al., 2022) and more recently for generative AI (Longpre et al., 2024). Additional proposals to facilitate external generative AI research include data donations (Sanderson, 2024).

Regulations like the Digital Services Act prescribe transparency by requiring audits of social media companies (European Commission, 2023), and similar auditing efforts are imaginable for generative AI. In fact, some scholars suggest to extend and adapt DSA rules for social media platforms to generative AI (Hacker et al., 2023).

While the specific implementation of these transparency efforts may be contentious and requires nuance, there is a broader lesson: Generative AI regulation can incentivize measures for increasing transparency, such as short and accessible explanations of the technology, independent oversight mechanisms, researcher access and mandatory audits.

#### 3.1.2. Gather Democratic Input to Inform Technology

Generative AI and social media share features that make them complex, including that the content they feature can pertain to a variety of domains, that there is potential for personalization, and that content could be moderated in various different ways. One approach to determine what a good content moderation system may look like is to gather input directly from users to inform design choices. Different initiatives have been launched over the past few years to gather input from users and enable democratic decisions about the nature of regulation and content moderation, with users deliberating issues ranging from cyberbullying on social platforms to the rules and constitutions that inform generative AI models (Wetherall-Grujić, 2023). These initiatives have their roots in the idea of deliberative democracy (Eagan, 2016). Social media also offers case studies of networks where content moderation seems to be broadly accepted and deliver productive results, such as in the case of the deliberation platform vTaiwan (Miller, 2019) or a neighborhood-focused social network (Oremus, 2024). Finally, researchers have studied how to embed important societal values into social media AI (Bernstein et al., 2023), which could inform how such values can be embedded into generative AI.

#### 3.1.3. PROMOTE USER CHOICE

Another option to empower users to make choices in the face of features such as content moderation and the varied nature of content is to enable users to set up rules for a subset of the system. The social media platform Mastodon is a prominent example in terms of increasing user choice in such a way. Mastodon is built on the idea that different communities can create their own servers and set and enforce their own content moderation rules (Mastodon, 2024). This highlights that the feature of personalization may be a potential route for resolving content moderation dilemmas. Content moderation questions with regards to generative AI and social media are similar and it is not clear what opinion representation should be the default, but increased personalization of models may be an answer (Redpoint, 2020).

# **3.2.** Address Specific Regulatory Concerns and Invest in Trust and Safety

The feature of content moderation that generative AI and social media share comes with challenges such as preventing the spread of harmful misinformation and protecting user wellbeing. Social media companies have invested in teams that address these specific regulatory concerns. Examples include teams at companies like Google, Meta and Microsoft working on youth wellbeing and mental health in general, election integrity, preventing spam, preventing the spread of child sexual abuse material, preventing harmful misinformation, detecting deceptive campaigns, and ensuring trust in the platform and safety of its users in general.

Generative AI chatbot performance has already been rated with regards to certain principles that apply just as much to social media. Common Sense Media published rankings of different generative AI models with respect to the following principles: put people first, prioritize fairness, be trustworthy, keep kids and teens safe, be effective, help people connect, use data responsibly, and be transparent and accountable (Common Sense Media, 2024).

Yet, generative AI companies do not have teams at the same scale as social media companies to address these issues. Generative AI companies are much smaller and younger than some of the social media giants, thus it is not surprising that they do not have as much dedicated staff to work on these issues. Going forward, however, adding diverse staff beyond engineers that can bring in expertise to address issues such as user mental health or combating misinformation is important to address the variety of risks and harms that generative AI models pose (for taxonomies of risks and harms related to generative AI, see Weidinger et al., 2021; Marchal et al., 2024; Gabriel et al., 2024; MITRE Corporation, 2021). Investment in trust and safety teams seems particularly crucial, and it is encouraging to see that companies like OpenAI and Anthropic are investing in this area, with OpenAI publishing the first-ever report on the activity of deceptive campaigns on generative AI platforms in May 2024 (Nimmo, 2024).

The policies social media companies have put in place to decide how and when to moderate individual users, and the best practices they have developed to uncover abuse such as deceptive campaigns that try to interfere with elections or spam users, could inform the approaches generative AI companies take. This includes developing a repertoire of content moderation approaches, which could include bans, but also more cautious interventions such as warnings and strikes for misbehavior, putting more guardrails in place or throttling usage for users that try to abuse generative AI models. Social media companies also gained experience in involving the user community in content moderation decisions (e.g., in the case of BirdWatch (Wojcik et al., 2022)) and how to collaborate across platforms, and generative AI companies could consider how these approaches could be adapted to their platforms.

Importantly, implementation of trust and safety measures for generative AI does not have to start from scratch. Opensource, collaborative tools like the Robust Open Online Safety Tools (ROOST, 2025) are a concrete example of collaboration across platforms and enable access to trust and safety resources even for companies with limited resources.

#### 3.3. Promote Computational Social Science Research

Both generative AI and social media allow users to express themselves and allow for a connection, be it to other users or to an AI with a vast pool of knowledge. How these media interact with users is a key part of what makes them so influential. They are neither purely technical, nor purely social systems. This suggests that multidisciplinary study computational social science—is needed to understand, evaluate and shape these systems (Gillespie et al., 2024).

In fact, the recommendations above, whether regarding measures to reduce bias or to enhance user wellbeing, all require computational social science research to test their effectiveness. Social media companies have hired researchers from many disciplines, including computer science, psychology, political science, communication, law and others, to better understand how their platforms impact society, and how certain interventions influence society and their revenue.

Rigorous computational social science evaluations, whether conducted in-house or via external researchers with platform access, are key to ensuring that technologies such as generative AI and social media meet their goal of being helpful and not harmful to society. Further investment in research is needed because generative AI has features that differ from previous technologies, so its impact and user preferences (e.g., with regards to privacy, personalization or content moderation) are not clear. Even the impact of previous technologies such as social media has not yet been comprehensively evaluated and needs further investment. Rigorous research can inform platform and public policy when it comes to regulation, and it can enhance user trust.

This implies the need to invest in diverse research teams that understand the interaction of humans and technology and can evaluate the societal implications technology. While AI company recruiting often focuses heavily on engineers, and some companies are more concerned with extreme risks in the more distant future, social media companies have shown the value of creating multidisciplinary teams to address current risks such as biases. Multidisciplinary teams allow companies to test different product features and interventions effectively, e.g. to reduce spam or misinformation spread. Guidance on building effective red teams for generative AI models also highlights the importance of diverse teams (Ofcom, 2024; Metcalf & Singh, 2024; Ahmad et al., 2024; Oremus, 2023). Computational social scientists from any background, data scientists and user experience researchers would be especially helpful to address questions at the intersection of humans and technology, such as which emotional bonds may be formed between humans and AI, and what type of personalization should be implemented.

While content moderation on social media is far from a resolved issue, there is a large and growing body of academic literature addressing user preferences and content moderation approaches (e.g., Persily & Tucker, 2020; Appel et al., 2023; Kozyreva et al., 2024), which could inform content moderation for generative AI.

#### 3.4. Take on a More Global Perspective

As the features spatial separation, general content, and use of AI imply, both generative AI and social media can be used in a variety of contexts. Generative AI companies have grown rapidly and are serving users around the world, similar to social media companies. However, compared to social media companies, many generative AI companies are more heavily focused on the US, likely due to their headquarter location (with exceptions such as Google DeepMind in the UK, Mistral AI in France, and DeepSeek in China). To address problems like biases, it is crucial that even small companies take on a global perspective and embrace local expertise in multiple countries. The reasoning mirrors that for the benefits of diversity in AI red teaming (Ofcom, 2024; Metcalf & Singh, 2024; Oremus, 2023), i.e. that broader representation allows for a better understanding of user preferences and the harms that a technology may pose. Taking on a more global perspective could take the form of establishing local offices and a focus on hiring internationally. The stakes are high. If companies fail to invest in taking user preferences and risk factors outside of the US seriously, the technology may serve large numbers of users worse (e.g., due to under-investment in non-English language content generation), contain undiscovered harms (Metcalf & Singh, 2024; Oremus, 2023), and could even result in catastrophes such as promoting violence in conflict regions (Amnesty International, 2022). Given the increasing amount of national and local regulations on generative AI, global expertise is also important to keep up with local laws.

For effective regulation, local expertise needs to be integrated into a global perspective. For example, the former Prime Minister of New Zealand suggested that a model for governing AI could follow the Christchurch Call, which is a multinational, multi-stakeholder effort bringing together governments, tech companies and civil society to eliminate violent extremist and terrorist content online (Ardern, 2023).

#### 4. Discussion

Applying learnings from social media regulation to generative AI regulation is challenging, in particular in light of political polarization, the rapid pace of technological development, the need to take different stakeholder characteristics into account, the patchwork of emerging AI regulations around the globe, and the need for effective implementation.

Policy is shaped by the contemporary political context. Thus, it is important to acknowledge that high levels of political polarization (documented e.g. in Finkel et al., 2020; Ruggeri et al., 2021) make it more difficult to develop policy that enjoys broad support (Druckman et al., 2021). When it comes to AI, public opinion and trust in governments and AI companies vary greatly based on people's characteristics, such as partisanship and nationality (Ipsos, 2024; Mcclain et al., 2025; Dreksler et al., 2025). Generative AI regulation has to be designed with this political context in mind. While polarization is a deeply-rooted issue that is difficult to address, our hope is that the recommendations on countering bias and perceptions thereof, increasing transparency and researcher access, and gathering democratic input to inform technology design may help depolarize the debate.

Another challenge for generative AI regulation is the pace of technology development. The capability leaps and proliferation of new models suggest that generative AI development outpaces social media development. This points to the need for more flexible policy, which is designed with foresight and is adaptable to future changes in the technology stack.

Further, effective regulation needs to take different stakeholder characteristics into account. For example, AI developers differ in their size, products, popularity, and resources. This implies that different stakeholders face different challenges in implementing regulations and managing compliance burden. Thus, regulations should be flexible. For example, the EU AI Act (European Parliament and Council, 2024) has different requirements depending on company size, user base size, compute used to create AI systems and whether AI systems are open source.

As argued in the previous section, generative AI regulation needs a global perspective. Yet, it has to account for a patchwork of different AI regulations around the globe (Alanoca et al., 2025). Some of these regulations are more compatible with the recommendations in this paper than others, which makes the implementation of these recommendations more feasible in some jurisdictions. The EU AI Act (European Parliament and Council, 2024) is one of the most comprehensive AI regulations around the globe. It regulates based on the level of risk that an AI system

poses and takes other factors such as openness of the technology into account. The EU AI Act aligns with several of the recommendations above, including transparency requirements, researcher access requirements, and a focus on specific regulatory concerns and risk areas. Work on the EU AI Act began before generative AI was widely adopted, but expanded its scope to include generative AI as it emerged. Thus, the EU AI Act addresses both traditional AI, such as social media, and generative AI, which may allow to integrate learnings from social media. Other regulation that was originally designed with technologies such as social media in mind is the EU Digital Services Act (DSA) (European Parliament and Council, 2022). This regulation also features transparency requirements, researcher access requirements, and specific guidelines when it comes to regulatory concerns such as election interference. Some argue that DSA rules could be adapted for generative AI platforms (Hacker et al., 2023). Another pertinent and controversial regulation that shaped social media is Section 230 of the Communications Decency Act (U.S. Congress, 1996). Section 230 holds that interactive computer service providers, including social media companies, are not considered publishers or speakers when they provide information that was provided to them by other users, which greatly limits platforms' liability for problematic content shared by users. This relates to generative AI regulation because AI developers face similar liability questions: While AI developers provide information that their technology generated and they could be considered publishers or creators, their content generation is based on training on others' speech and they could be seen as intermediaries. Section 230 has been the subject of intense legal and political debate, with uncertain outcomes. The debates about Section 230 demonstrate that definitions and liability determinations can evolve and can shape platforms' business models and content moderation strategies. Learning from existing regulation before adding new ones can help address the lack of specificity of the current environment, and prevent further fragmentation.

A final challenge for effective regulation is effective implementation. Regulation should be concrete enough to allow for smooth implementation by stakeholders with different characteristics—regulations and standards that are too vague impose a significant burden on stakeholders and are difficult to monitor (Pouget & Zuhdi, 2024). More generally, it is crucial to track readiness and compliance (see e.g. Scott, 2024; Bommasani et al., 2023b) and ensure proper evaluation and incentives.

#### **5.** Alternative Views

In this paper, we argued that generative AI regulation can learn from social media regulation. However, there are valid counterarguments related to the imperfect analogy between generative AI and social media, the fact that social media regulation has not been a model example of technology regulation, and that thinking about regulation from first principles may be desirable.

First, as discussed in detail in Section 2.2, there are important differences in the affordances of social media and generative AI, including whether the medium acts as conversation partner, how visible content moderation is, and how many functions the medium serves. For example, generative AI and social media differ in that only generative AI chatbots tend to be a direct conversation partner for humans, in contrast to traditional social media platforms where people post to interact with other people. This could have implications for the kind of relationships that people form with the technology, which in turn could affect the need for regulation. While this concern is valid, the recommendations in this paper build on the affordances that generative AI and social media share and do not go beyond those. For differences in affordances, other comparisons could be insightful (for a review of different AI metaphors, see Maas, 2023). For example, learning from regulatory authorities such as the FDA could help inform AI governance (Raji et al., 2022; AI Now Institute, 2024).

Second, as described in earlier parts such as Section 3.1.1, social media regulation has not been a model example of technology regulation. Research projects like the Facebook and Instagram Election Study (Clegg & Nayak, 2020) faced major delays, and initiatives like alternative oversight boards (The Real Facebook Oversight Board, 2022) and coalitions to protect independent researchers (Coalition for Independent Technology Research, 2022) show that the research community and the broader public have not been satisfied with how social media regulation played out. However, we can learn lessons from both past failures and past successes. The encouragement to learn from social media regulation does not mean that we should always take similar regulatory approaches for generative AI. It means that we should carefully assess what worked well, and what needs to be improved, to let these insights inform generative AI regulation.

Third, instead of looking into the past, it may be desirable to think about generative AI regulation from first principles (Clear, 2024). This could prevent getting caught up in unhelpful norms and precedents that may prevent innovative and effective regulation. This is a valid point, but learning from social media regulation and thinking about generative AI regulation from first principles are not mutually exclusive. For example, it is possible to start thinking about desirable regulation free from any other existing ideas, and afterwards analyze whether the developed approaches are promising in light of what we know about technology regulation in other areas. Learning lessons from social media is the best way to prevent avoidable mistakes because many challenges that generative AI regulation aims to address, including issues such as content moderation and bias, are not as unprecedented as they may seem.

## 6. Conclusion

There are strong disagreements about the approach that should be taken to regulate generative AI. This paper argued that the regulation of generative AI can be informed by the evolution of the regulation of social media. While social media is not the only analogy proposed for generative AI (Maas, 2023), and by no means a perfect analogy, generative AI and social media share key features that make a comparison of the two worthwhile. An analysis of social media regulation efforts-including self-regulation and laws-reveals interesting approaches and best practices. This paper outlined recommendations regarding transparency, researcher access, gathering democratic input, promoting user choice, addressing specific regulatory concerns, increasing investments into computational social science research, and taking on a more global perspective. In the case of social media, self-regulation did not always work, which has resulted in multiple new laws being proposed in the past few years. These self-regulation efforts and laws, including specific approaches to increasing transparency, enhancing user choice, and investing in research, can be valuable pointers for those looking to regulate generative AI. Analyzing social media regulation may inform and accelerate the process of developing generative AI regulation. Regulation takes time and effort, so where possible, we should save resources and avoid mistakes by learning the lessons that social media regulation holds for generative AI regulation.

#### Acknowledgements

The author is grateful to Jennifer Pan for feedback on an earlier draft.

#### **Impact Statement**

Against the backdrop of increasingly heated debates about generative AI regulation, this paper shows that we do not have to reinvent the wheel when it comes to questions such as how to ensure that generative AI is safe and moderated in alignment with users' preferences. Instead, we can learn lessons from social media regulation.

Concrete lessons we can learn include the importance of investing in trust and safety and taking a more diverse perspective, both in terms of geography and research disciplines.

Learning lessons from social media regulation can help prevent avoidable mistakes and utilize resources more effectively, which can ultimately improve AI policy and AI safety.

#### References

- Abdo, A., Krishnan, R., Krent, S., and Woods, A. K. A safe harbor for platform research, 2022. URL https://knightcolumbia.org/content/asafe-harbor-for-platform-research.
- Abid, A., Farooqi, M., and Zou, J. Large language models associate Muslims with violence. *Nature Machine Intelligence*, 3(6):461–463, 2021. ISSN 25225839. doi: 10.1038/s42256-021-00359-2. URL https://doi.org/10.1038/s42256-021-00359-2.
- Ahmad, L., Agarwal, S., Lampe, M., and Mishkin, P. OpenAI's approach to external red teaming for AI models and systems. Technical report, OpenAI, 1 2024. URL https://cdn.openai.com/ papers/openais-approach-to-externalred-teaming.pdf.
- AI Now Institute. Lessons from the FDA for AI. Technical Report August, AI Now Institute, 2024. URL https://ainowinstitute.org/wpcontent/uploads/2024/08/20240801-AI-Now-FDA.pdf.
- Alanoca, S., Gur-Arieh, S., Zick, T., and Klyman, K. Comparing apples to oranges: A taxonomy for navigating the global landscape of AI regulation, 2025. URL https: //doi.org/10.48550/arXiv.2505.13673.
- Amnesty International. The social atrocity: Meta and the right to remedy for the Rohingya. Technical report, Amnesty International, 2022. URL https://www.amnesty.org/en/documents/ asa16/5933/2022/en/.
- Anthropic. Measuring the persuasiveness of language models, 2024. URL https://www.anthropic.com/ news/measuring-model-persuasiveness.
- Appel, R. E., Pan, J., and Roberts, M. E. Partisan conflict over content moderation is more than disagreement about facts. *Science Advances*, 9(44):1–10, 2023. URL https: //doi.org/10.1126/sciadv.adg6799.
- Ardern, J. There's a model for governing AI. Here it is., jun 2023. URL https://www.washingtonpost. com/opinions/2023/06/09/jacindaardern-ai-new-zealand-planning/.
- Bai, H., Voelkel, J. G., Eichstaedt, J. C., and Willer, R. Artificial intelligence can persuade humans on political issues, Feb 2023. URL osf.io/stakv\_v1.
- Barrett, P. M. and Sims, J. G. False accusation: The unfounded claim that social media companies censor conservatives. Technical Report February,

NYU Stern Center for Business and Human Rights, 2021. URL https://static1.squarespace. com/static/5b6df958f8370af3217d4178/ t/6011e68dec2c7013d3caf3cb/ 1611785871154/NYU+False+Accusation+ report\_FINAL.pdf.

- Bereska, L. and Gavves, E. Mechanistic interpretability for ai safety – a review, 2024. URL https://arxiv. org/abs/2404.14082.
- Bernstein, M. S., Christin, A., Hancock, J. T., Hashimoto, T., Jia, C., Lam, M., Persily, N., Piccardi, T., Saveski, M., Tsai, J. L., Ugander, J., and Xu, C. Embedding societal values into social media algorithms. *Journal of Online Trust and Safety*, 2(1):1–13, 2023. URL https: //doi.org/10.54501/jots.v2i1.148.
- Bommasani, R., Klyman, K., Longpre, S., Kapoor, S., Maslej, N., Xiong, B., Zhang, D., and Liang, P. The Foundation Model Transparency Index, 2023a. URL http://arxiv.org/abs/2310.12941.
- Bommasani, R., Klyman, K., Zhang, D., and Liang, P. Do foundation model providers comply with the EU AI Act?, 2023b. URL https://crfm.stanford. edu/2023/06/15/eu-ai-act.html.
- Bommasani, R., Klyman, K., Kapoor, S., Longpre, S., Xiong, B., Maslej, N., and Liang, P. The Foundation Model Transparency Index v1.1, 2024. URL https: //doi.org/10.48550/arXiv.2407.12929.
- Bostrom, N. Existential risk prevention as global priority. *Global Policy*, 4(1):15–31, 2013. doi: https://doi.org/10.1111/1758-5899.12002. URL https://onlinelibrary.wiley.com/doi/ abs/10.1111/1758-5899.12002.
- Cattell, S., Ghosh, A., and Kaffee, L.-A. Coordinated disclosure for AI: Beyond security vulnerabilities, 2024. URL http://arxiv.org/abs/2402.07039.
- Center for Humane Technology. The attention economy: Why do tech companies fight for our attention? Technical report, Center for Humane Technology, 2021. URL https://cdn.prod.websitefiles.com/5f0e1294f002b15080e1f2ff/ 612f8e3fa20df8374659a774\_1-TheAttentionEconomyIssueGuide.pdf.
- Clark, H. H. Using Language. Cambridge University Press, 1996. URL https://doi.org/10.1017/ CB09780511620539.
- Clear, J. First principles: Elon Musk on the power of thinking for yourself, 2024. URL https://jamesclear. com/first-principles.

- Clegg, N. How AI influences what you see on Facebook and Instagram, 2023. URL https://about.fb.com/ news/2023/06/how-ai-ranks-content-onfacebook-and-instagram/.
- Clegg, N. and Nayak, C. New Facebook and Instagram research initiative to look at US 2020 presidential election, August 2020. URL https: //about.fb.com/news/2020/08/researchimpact-of-facebook-and-instagram-onus-election/.
- Coalition for Independent Technology Research. Coalition for Independent Technology Research Founding Document, 2022. URL https://independenttechresearch.org/ coalition-for-independent-technologyresearch-founding-document/.
- Common Sense Media. AI Initiative, 2024. URL https: //www.commonsensemedia.org/ai.
- Dominguez-Olmedo, R., Hardt, M., and Mendler-Dünner, C. Questioning the survey responses of large language models, 2024. URL https://arxiv.org/abs/2306. 07951.
- Dreksler, N., Law, H., Ahn, C., Schiff, D., Schiff, K., and Peskowitz, Z. What Does the Public Think About AI? Technical report, Centre for the Governance of AI, 2025. URL https://doi.org/10.2139/ssrn. 5108572.
- Druckman, J. N., Klar, S., Krupnikov, Y., Levendusky, M., and Ryan, J. B. Affective polarization, local contexts and public opinion in America. *Nature Human Behaviour*, 5 (1):28–38, January 2021. URL https://doi.org/ 10.1038/s41562-020-01012-5.
- Durmus, E., Nyugen, K., Liao, T. I., Schiefer, N., Askell, A., Bakhtin, A., Chen, C., Hatfield-Dodds, Z., Hernandez, D., Joseph, N., Lovitt, L., McCandlish, S., Sikder, O., Tamkin, A., Thamkul, J., Kaplan, J., Clark, J., and Ganguli, D. Towards measuring the representation of subjective global opinions in language models, 2023. URL http://arxiv.org/abs/2306.16388.
- Eagan, J. L. deliberative democracy, 2016. URL https://www.britannica.com/topic/ deliberative-democracy.
- European Commission. Shaping Europe's digital future Commission adopts rules on independent audits under the Digital Services Act, 2023. URL https://digital-strategy.ec.europa. eu/en/news/commission-adopts-rulesindependent-audits-under-digitalservices-act.

- European Parliament and Council. Regulation (EU) 2022/2065 of the European Parliament and of the Council of 19 October 2022 on a Single Market For Digital Services and amending Directive 2000/31/EC (Digital Services Act), 2022. URL https://eur-lex.europa.eu/eli/reg/2022/2065/oj/eng.
- European Parliament and Council. Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828 (Artificial Intelligence Act), 2024. URL https://eur-lex.europa.eu/legal-content/EN/TXT/HTML/?uri=OJ: L\_202401689.
- Executive Office of the President. Safe, secure, and trustworthy development and use of artificial intelligence. Executive Order, 10 2023. URL https://www.federalregister.gov/ documents/2023/10/30/2023-24110/safesecure-and-trustworthy-developmentand-use-of-artificial-intelligence. Federal Register Vol. 88, No. 210 (October 30, 2023).
- Executive Office of the President. Removing barriers to american leadership in artificial intelligence. Executive Order, 1 2025. URL https://federalregister.gov/d/2025-02172.
- Ferrara, E. Should ChatGPT be biased? challenges and risks of bias in large language models, 2023. URL http://arxiv.org/abs/2304.03738.
- Finkel, E. J., Bail, C. A., Cikara, M., Ditto, P. H., Iyengar, S., Klar, S., Mason, L., McGrath, M. C., Nyhan, B., Rand, D. G., Skitka, L. J., Tucker, J. A., Bavel, J. J. V., Wang, C. S., and Druckman, J. N. Political sectarianism in america. *Science*, 370(6516):533–536, 2020. doi: 10.1126/ science.abe1715. URL https://www.science. org/doi/abs/10.1126/science.abe1715.
- Fisher, J., Feng, S., Aron, R., Richardson, T., Choi, Y., Fisher, D. W., Pan, J., Tsvetkov, Y., and Reinecke, K. Biased AI can influence political decision-making, 2024. URL https://arxiv.org/abs/2410.06415.
- Gabriel, I., Manzini, A., Keeling, G., Hendricks, L. A., Rieser, V., Iqbal, H., Tomašev, N., Ktena, I., Kenton, Z., Rodriguez, M., El-Sayed, S., Brown, S., Akbulut, C., Trask, A., Hughes, E., Bergman, A. S., Shelby, R., Marchal, N., Griffin, C., Mateos-Garcia, J., Weidinger, L., Street, W., Lange, B., Ingerman, A., Lentz, A., Enger, R., Barakat, A., Krakovna, V., Siy, J. O.,

Kurth-Nelson, Z., Mccroskery, A., Bolina, V., Law, H., Shanahan, M., Alberts, L., Balle, B., De Haas, S., Ibitoye, Y., Dafoe, A., Goldberg, B., Krier, S., Reese, A., Witherspoon, S., Hawkins, W., Rauh, M., Wallace, D., Franklin, M., Goldstein, J. A., Lehman, J., Klenk, M., Vallor, S., Biles, C., Morris, M. R., King, H., Agüera Y Arcas, B., Isaac, W., and Manyika, J. The ethics of advanced AI assistants, 2024. URL https: //storage.googleapis.com/deepmindmedia/DeepMind.com/Blog/ethics-ofadvanced-ai-assistants/the-ethics-ofadvanced-ai-assistants-2024-i.pdf.

- Gillespie, T., Shaw, R., Gray, M. L., and Suh, J. AI redteaming is a sociotechnical system. now what?, 2024. URL https://arxiv.org/abs/2412.09751.
- Hacker, P., Engel, A., and Mauer, M. Regulating Chat-GPT and other large generative AI models. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '23, pp. 1112–1123, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701924. doi: 10.1145/ 3593013.3594067. URL https://doi.org/10. 1145/3593013.3594067.
- Hendrycks, D., Schmidt, E., and Wang, A. Superintelligence strategy: Expert version, 2025. URL https://arxiv. org/abs/2503.05628.
- Hu, T., Kyrychenko, Y., Rathje, S., Collier, N., van der Linden, S., and Roozenbeek, J. Generative language models exhibit social identity biases. *Nature Computational Science*, 2024. ISSN 26628457. doi: 10.1038/s43588-024-00741-1. URL https://doi.org/10.1038/ s43588-024-00741-1.
- Huttenlocher, D., Ozdaglar, A., and Goldston, D. A framework for U.S. AI governance: Creating a safe and thriving AI sector. Technical report, MIT Schwarzman College of Computing, 2023. URL https://computing.mit.edu/wp-content/ uploads/2023/11/AIPolicyBrief.pdf.
- Ipsos. The Ipsos AI Monitor: A 32-country Ipsos Global Advisor Survey. Technical Report June, Ipsos, 2024. URL https://www.ipsos.com/sites/default/ files/ct/news/documents/2024-06/ Ipsos-AI-Monitor-2024-final-APAC.pdf.
- Isaac, M. and Schleifer, T. Meta says it will end its fact-checking program on social media posts. The New York Times, January 2025. URL https://www.nytimes.com/live/2025/ 01/07/business/meta-fact-checking. Updated January 15, 2025.

- Iyer, R. To evaluate meta's shift, focus on the product changes, not the moderation. *Tech Policy Press*, January 2025. URL https://www.techpolicy.press/ to-evaluate-metas-shift-focus-on-theproduct-changes-not-the-moderation/.
- Jakesch, M., Bhat, A., Buschek, D., Zalmanson, L., and Naaman, M. Co-writing with opinionated language models affects users' views. In *Conference on Human Factors in Computing Systems*, 2023. doi: 10. 1145/3544548.3581196. URL https://doi.org/ 10.1145/3544548.3581196.
- Kozyreva, A., Lorenz-Spreen, P., Herzog, S. M., Ecker, U. K., Lewandowsky, S., Hertwig, R., Ali, A., Bak-Coleman, J., Barzilai, S., Basol, M., Berinsky, A. J., Betsch, C., Cook, J., Fazio, L. K., Geers, M., Guess, A. M., Huang, H., Larreguy, H., Maertens, R., Panizza, F., Pennycook, G., Rand, D. G., Rathje, S., Reifler, J., Schmid, P., Smith, M., Swire-Thompson, B., Szewach, P., van der Linden, S., and Wineburg, S. Toolbox of individual-level interventions against online misinformation. *Nature Human Behaviour*, pp. 1–9, 2024. ISSN 23973374. URL https://doi.org/10.1038/s41562-024-01881-0.
- Kumar, V. Making "freemium" work. *Harvard Business Review*, May 2014. URL https://hbr.org/2014/ 05/making-freemium-work.
- Li, D., Pyke, R., Jiang, R., and Jagadeesh, K. Introducing the Researcher Platform: Empowering independent research analyzing large-scale data from Meta, January 2022. URL https://research. facebook.com/blog/2022/1/introducingthe-researcher-platform-empoweringindependent-research-analyzing-largescale-data-from-meta/.
- Longpre, S., Kapoor, S., Klyman, K., Ramaswami, A., Bommasani, R., Blili-Hamelin, B., Huang, Y., Skowron, A., Yong, Z. X., Kotha, S., Zeng, Y., Shi, W., Yang, X., Southen, R., Robey, A., Chao, P., Yang, D., Jia, R., Kang, D., Pentland, A., Narayanan, A., Liang, P., and Henderson, P. Position: A safe harbor for AI evaluation and red teaming. In Salakhutdinov, R., Kolter, Z., Heller, K., Weller, A., Oliver, N., Scarlett, J., and Berkenkamp, F. (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 32691–32710. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/longpre24a.html.
- Lukoff, K., Lyngs, U., Zade, H., Liao, J. V., Choi, J., Fan, K., Munson, S. A., and Hiniker, A. How the design of YouTube influences user sense of agency.

In Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems, CHI '21, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450380966. doi: 10.1145/ 3411764.3445467. URL https://doi.org/10. 1145/3411764.3445467.

- Maas, M. M. AI is like... a literature review of AI metaphors and why they matter for policy. *SSRN Electronic Journal*, October 2023. doi: 10.2139/ssrn.4612468. URL https://doi.org/10.2139/ssrn.4612468.
- Marchal, N., Xu, R., Elasmar, R., Gabriel, I., Goldberg, B., and Isaac, W. Generative AI misuse: A taxonomy of tactics and insights from real-world data, 2024. URL https://arxiv.org/abs/2406.13843.
- Mastodon. Mastodon, 2024. URL https://joinmastodon.org/.
- Matz, S. C., Teeny, J., Vaid, S., Harari, G., and Cerf, M. The potential of generative AI for personalized persuasion at scale. *Nature Scientific Reports*, 14(4692): 1–16, 2024. URL https://doi.org/10.1038/ s41598-024-53755-0.
- Mcclain, B. Y. C., Kennedy, B., Gottfried, J., and Anderson, M. How the U.S. public and AI experts view artificial intelligence. Technical Report April, Pew Research Center, 2025. URL https://www.pewresearch.org/wpcontent/uploads/sites/20/2025/04/ pi\_2025.04.03\_us-public-and-aiexperts\_report.pdf.
- McGrenere, J. and Ho, W. Affordances: Clarifying and evolving a concept. *Graphics Interface*, pp. 1–8, 2000. URL https://graphicsinterface.org/ wp-content/uploads/gi2000-24.pdf.
- Meta. Our approach to Facebook feed ranking how feed ranking works for connected content, 2024a. URL https://transparency.fb.com/ features/ranking-and-content/.
- Meta. Oversight Board, 2024b. URL https:// transparency.fb.com/en-gb/oversight/ oversight-board-recommendations/.
- Metcalf, J. and Singh, R. Scaling up mischief: Redteaming AI and distributing governance. *Harvard Data Science Review*, May 2024. URL https://hdsr. mitpress.mit.edu/pub/ded4vcwl.
- Miller, C. Crossing Divides: How a social network could save democracy from deadlock, 2019. URL https:// www.bbc.com/news/technology-50127713.

- MITRE Corporation. MITRE ATLAS: Adversarial threat landscape for artificial-intelligence systems, 2021. URL https://atlas.mitre.org/.
- Nangia, N., Vania, C., Bhalerao, R., and Bowman, S. R. CrowS-Pairs: A challenge dataset for measuring social biases in masked language models, 2020. URL https://aclanthology.org/2020. emnlp-main.154/.
- Naous, T., Ryan, M. J., and Xu, W. Having beer after prayer? measuring cultural bias in large language models, 2023. URL https://arxiv.org/abs/2305. 14456v1.
- Nimmo, B. AI and covert influence operations: Latest trends. Technical Report May, OpenAI, 2024. URL https://downloads.ctfassets.net/ kftzwdyauwt9/5IMxzTmUclSOAcWUXbkVrK/ 3cfab518e6b10789ab8843bcca18b633/ Threat\_Intel\_Report.pdf.
- Ofcom. Red teaming for GenAI harms: Revealing the risks and rewards for online safety. Discussion paper, Office of Communications, 7 2024. URL https://www.ofcom.org.uk/siteassets/ resources/documents/consultations/ discussion-papers/red-teaming/redteaming-for-gen-ai-harms.pdf.
- Oremus, W. Meet the hackers who are trying to make AI go rogue, 2023. URL https: //www.washingtonpost.com/technology/ 2023/08/08/ai-red-team-defcon/.
- Oremus, W. The friendliest social network you've never heard of, August 2024. URL https://www.washingtonpost.com/ technology/2024/08/10/front-porchforum-vermont-research-new-public/.
- Persily, N. and Tucker, J. A. (eds.). Social Media and Democracy. Cambridge University Press, 2020. ISBN 9789162893439. URL https://www.cambridge.org/core/ books/social-media-and-democracy/ E79E2BBF03C18C3A56A5CC393698F117.
- Potter, Y., Lai, S., Kim, J., Evans, J., and Song, D. Hidden persuaders: LLMs' political leaning and their influence on voters. In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pp. 4244–4275, Miami, Florida, USA, November 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024. emnlp-main.244/.

- Pouget, H. and Zuhdi, R. AI and product safety standards under the EU AI Act. Carnegie Endowment for International Peace, March 2024. URL https://carnegieendowment.org/ research/2024/03/ai-and-productsafety-standards-under-the-eu-ai-act.
- Rafaeli, S. and Sudweeks, F. Networked interactivity. Journal of Computer-Mediated Communication, 2(4): JCMC243, March 1997. URL https://doi.org/ 10.1111/j.1083-6101.1997.tb00201.x.
- Raji, I. D., Xu, P., Honigsberg, C., and Ho, D. Outsider oversight: Designing a third party audit ecosystem for AI governance. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '22, pp. 557–571, New York, NY, USA, 2022. Association for Computing Machinery. URL https://doi.org/10. 1145/3514094.3534181.
- Redpoint. Stanford professor Tatsu Hashimoto on AI biases and improving LLM performance, 2020. URL https: //www.youtube.com/watch?v=pceYeZdT1D0.
- Robertson, A. Google apologizes for 'missing the mark' after Gemini generated racially diverse Nazis. *The Verge*, 2024. URL https://www.theverge.com/ 2024/2/21/24079371/google-ai-geminigenerative-inaccurate-historical.
- Romm, T. Senate Republicans renew their claims that Facebook, Google and Twitter censor conservatives, apr 2019. URL https://www.washingtonpost. com/technology/2019/04/10/facebookgoogle-twitter-under-fire-senaterepublicans-censoring-conservativesonline/.
- Ronzhyn, A., Cardenal, A. S., and Rubio, A. B. Defining affordances in social media research: A literature review. *New Media & Society*, 25(11):3165–3188, 2023. doi: 10.1177/14614448221135187. URL https://doi. org/10.1177/14614448221135187.
- ROOST. Robust open online safety tools, 2025. URL https://roost.tools/.
- Röttger, P., Hofmann, V., Pyatkin, V., Hinck, M., Kirk, H. R., Schütze, H., and Hovy, D. Political compass or spinning arrow? towards more meaningful evaluations for values and opinions in large language models, 2024. URL http://arxiv.org/abs/2402.16786.
- Rozado, D. The political biases of ChatGPT. Social Sciences, 12(3), 2023. URL https://doi.org/10.3390/socsci12030148.

- Ruggeri, K., Većkalov, B., Bojanić, L., Andersen, T. L., Ashcroft-Jones, S., Ayacaxli, N., Barea-Arroyo, P., Berge, M. L., Bjørndal, L. D., Bursalioğlu, A., Bühler, V., Čadek, M., Çetinçelik, M., Clay, G., Cortijos-Bernabeu, A., Damnjanović, K., Dugue, T. M., Esberg, M., Esteban-Serna, C., Felder, E. N., Friedemann, M., Frontera-Villanueva, D. I., Gale, P., Garcia-Garzon, E., Geiger, S. J., George, L., Girardello, A., Gracheva, A., Gracheva, A., Guillory, M., Hecht, M., Herte, K., Hubená, B., Ingalls, W., Jakob, L., Janssens, M., Jarke, H., Kácha, O., Kalinova, K. N., Karakasheva, R., Khorrami, P. R., Lep, Ž., Lins, S., Lofthus, I. S., Mamede, S., Mareva, S., Mascarenhas, M. F., McGill, L., Morales-Izquierdo, S., Moltrecht, B., Mueller, T. S., Musetti, M., Nelsson, J., Otto, T., Paul, A. F., Pavlović, I., Petrović, M. B., Popović, D., Prinz, G. M., Razum, J., Sakelariev, I., Samuels, V., Sanguino, I., Say, N., Schuck, J., Soysal, I., Todsen, A. L., Tünte, M. R., Vdovic, M., Vintr, J., Vovko, M., Vranka, M. A., Wagner, L., Wilkins, L., Willems, M., Wisdom, E., Yosifova, A., Zeng, S., Ahmed, M. A., Dwarkanath, T., Cikara, M., Lees, J., and Folke, T. The general fault in our fault lines. Nature Hu*man Behaviour*, 5(10):1369–1380, 2021. URL https: //doi.org/10.1038/s41562-021-01092-x.
- Sanderson, Z. Beyond competition: Designing data portability to support research on the digital information environment. *SSRN Electronic Journal*, 2024. URL https://doi.org/10.2139/ssrn.4739362.
- Santurkar, S., Durmus, E., Ladhak, F., Lee, C., Liang, P., and Hashimoto, T. Whose opinions do language models reflect? In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pp. 29971–30004. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/ v202/santurkar23a.html.
- Scott, M. 5 things to know about the Digital Services Act's first risk assessments and audits. TechPolicy.Press, December 2024. URL https: //www.techpolicy.press/5-things-toknow-about-the-digital-services-actsfirst-risk-assessments-and-audits/.
- Sharma, S. and Murano, P. A usability evaluation of web user interface scrolling types. *First Monday*, 25(3), February 2020. URL https://doi.org/10.5210/fm. v25i3.10309.
- Stern, J. AI is like ... nuclear weapons? The Atlantic, March 2023. URL https://www.theatlantic. com/technology/archive/2023/03/aigpt4-technology-analogy/673509/.

- Suresh, H. and Guttag, J. A framework for understanding sources of harm throughout the machine learning life cycle. In *Equity and Access in Algorithms, Mechanisms, and Optimization*, EAAMO '21, pp. 1–9. ACM, October 2021. doi: 10.1145/3465416.3483305. URL https: //doi.org/10.1145/3465416.3483305.
- The Real Facebook Oversight Board. The Real Facebook Oversight Board, 2022. URL https: //the-citizens.com/real-facebookoversight/.
- Thulasi, S. Understand why you're seeing certain ads and how you can adjust your ad experience, 2019. URL https://about.fb.com/news/2019/07/ understand-why-youre-seeing-ads/.
- TikTok. Research tools TikTok for developers, 2025.URL https://developers.tiktok.com/ products/research-api/.
- U.S. Congress. 47 U.S.C. § 230 Protection for private blocking and screening of offensive material. United States Code, 1996. URL http://uscode.house. gov/view.xhtml?req=(title:47+section: 230+edition:prelim).
- Wagner, M. W. Independence by permission. *Science*, 381 (6656):388–391, 2023. URL https://doi.org/10.1126/science.adi2430.
- Wartella, E. and Reeves, B. Historical trends in research on children and the media: 1900-1960. Journal of Communication, 35(2):118–133, 1985. URL https://doi.org/10.1111/j. 1460-2466.1985.tb02238.x.
- Weidinger, L., Mellor, J., Rauh, M., Griffin, C., Uesato, J., Huang, P.-S., Cheng, M., Glaese, M., Balle, B., Kasirzadeh, A., Kenton, Z., Brown, S., Hawkins, W., Stepleton, T., Biles, C., Birhane, A., Haas, J., Rimell, L., Hendricks, L. A., Isaac, W., Legassick, S., Irving, G., and Gabriel, I. Ethical and social risks of harm from language models, 2021. URL https://arxiv.org/ abs/2112.04359.
- Wetherall-Grujić, G. The race to democratise AI, 2023. URL https://democracytechnologies.org/participation/therace-to-democratise-ai/.
- Wiener, Roth, Rubio, and Stern. Safe and Secure Innovation for Frontier Artificial Intelligence Models Act, 2024. URL https://leginfo.legislature.ca. gov/faces/billTextClient.xhtml?bill\_ id=202320240SB1047.

- Williams-Ceci, S., Jakesch, M., Bhat, A., Kadoma, K., Zalmanson, L., and Naaman, M. Bias in AI autocomplete suggestions leads to attitude shift on societal issues, 2024. URL https://doi.org/10.31234/osf. io/mhjn6.
- Wojcik, S., Hilgard, S., Judd, N., Mocanu, D., Ragain, S., Hunzaker, M. B. F., Coleman, K., and Baxter, J. Birdwatch: Crowd wisdom and bridging algorithms can inform understanding and reduce the spread of misinformation, 2022. URL https://arxiv.org/abs/2210. 15723.
- Zaidan, E. and Ibrahim, I. A. AI governance in a complex and rapidly changing regulatory landscape: A global perspective. *Humanities and Social Sciences Communications*, 11(1):1–18, 2024. URL https://doi.org/ 10.1057/s41599-024-03560-x.