

# Research can help to tackle AI-generated disinformation

Stefan Feuerriegel, Renée DiResta, Josh A. Goldstein, Srijan Kumar, Philipp Lorenz-Spreen, Michael Tomz & Nicolas Pröllochs



Generative artificial intelligence (AI) tools have made it easy to create realistic disinformation that is hard to detect by humans and may undermine public trust. Some approaches used for assessing the reliability of online information may no longer work in the AI age. We offer suggestions for how research can help to tackle the threats of AI-generated disinformation.

In March 2023, images of former president Donald Trump ostensibly getting arrested circulated on social media. Former president Trump, however, did not get arrested in March. The [images were fabricated](#) using generative AI technology. Although the phenomenon of fabricated or altered content is not new, recent advances in generative AI technology have made it easy to produce fabricated content that is increasingly realistic, which makes it harder for people to distinguish what is real.

Generative AI tools can be used to create original content, such as text, images, audio and video. Although most applications of these tools are benign, there is substantial concern about the potential for increased proliferation of disinformation (which we refer to broadly as content spread with the intent to deceive, including propaganda and fake news). Because the content generated appears highly realistic, some of the strategies presently used for detecting manipulative accounts and content are rendered ineffective by AI-generated disinformation.

## How AI disinformation differs

What makes AI-generated disinformation different from traditional, human-generated disinformation? Here, we highlight four potentially differentiating factors: scale, speed, ease of use and personalization. First, generative AI tools make it possible to mass-produce content for disinformation campaigns. One example of the scale of AI-generated disinformation is the use of generative AI tools to produce [dozens of different fake images](#) showing Pope Francis in haute fashion across different postures and backgrounds. In particular, AI tools can be used to create multiple variations of the same false stories, translate them into different languages, mimic conversational dialogues and more. Second, compared to the manual generation of content, AI technology allows disinformation to be produced very rapidly. For example, fake images [can be created with tools such as Midjourney](#) in seconds, whereas without generative AI the creation of similar images would take hours or days. These first two factors – scale and speed – are

challenges for fact-checkers, who will be flooded with disinformation but still need substantial amounts of time for debunking. Third, as AI tools diffuse into society more broadly, they will lower the barriers to entry for running influence operations. People can use AI tools to create realistic fake images and videos without professional expertise or time-consuming manual editing. This may democratize the troll farm. Fourth, AI technology renders it easier to launch personalized disinformation campaigns to target-specific audiences (or individuals) and their preferences or beliefs without deep knowledge of the language or culture of the target. For example, personalized disinformation may target people of different ages, political ideologies, religious beliefs and personality types (for example, such as extroverts or introverts), which may increase the persuasiveness of disinformation campaigns. Those already marginalized by society or who have low media literacy may be particularly vulnerable.

The scalability of AI technology could enhance the tactics of disinformation campaigns. First, tactics that involve highly targeted one-to-one communication (for example, through bots or other automated tools) may become more common. For example, scammers may create generated audio content that [sounds like a distressed family member](#) to demand ransoms from targeted victims. Second, the scale of content production may augment tactics aimed at distracting audiences and at creating the illusion of majorities (as content appears to be coming from different sources). For example, state and state-linked actors (such as Russia's Internet Research Agency) have long leveraged hundreds of accounts to [divert attention from inconvenient stories](#), which will now be easier with generative AI. Third, creating back-and-forth conversations in real time may help to obscure the automated nature of corresponding social media accounts.

It is important to note that the relevant threat vectors are broader than social media: AI tools enable low-cost and high-volume fabrication of email campaigns<sup>1</sup>, paid advertisements, websites and scientific documents that provide false evidence for claims. At the extreme, the deluge of AI-generated disinformation may make it more difficult to discern the truth from the noise in online spaces, and reduce broader societal trust.

## AI disinformation is hard to detect

Existing research shows that generative AI systems can write disinformation that is hard to detect by humans because it can mimic the style of reliable sources and communicators who are trusted by the target audience<sup>2,3</sup>. Research is needed to explore the vulnerabilities of individuals to AI-generated disinformation campaigns and, specifically, the extent to which capabilities such as microtargeting and one-on-one chats make disinformation campaigns more credible and persuasive. For example, future research could compare differences in behavioural outcomes between persuasive and distractive campaigns with AI-generated disinformation.

Endogenous cues<sup>4</sup> that humans have previously used for judging the reliability of information (for example, whether a text is free of grammatical errors or whether images have lighting and shadows that are consistent with reality) may no longer be good determinants. As a result, the importance of non-content-based, exogenous cues may increase. These cues include indications of who wrote or created content, whether it bears a watermark, how well a piece of information is connected to existing knowledge (for example, by references to known sources)<sup>4</sup> or how other people in one's network interact with it (for example, diversity of readership is associated with quality)<sup>5</sup>. We believe that one high-value research direction is identifying and testing exogenous cues that are transparent and difficult to game with generative AI.

## Behavioural mitigations

To mitigate the effects of disinformation campaigns, behavioural interventions can be used that involve nudging the behaviour of online users or boosting their competencies so they are less likely to believe or share fake content. However, it is not clear how effective previously designed behavioural interventions will be in tackling AI-generated disinformation. For the reasons discussed above, behavioural interventions based on endogenous cues (such as asking people to check whether images look realistic before sharing them or asking people to check for professional writing quality) may be insufficient, as generative AI can fabricate images that look authentic and text that is largely error-free. Rather, behavioural interventions aimed at exogenous cues will grow more important. We see at least three categories of behavioural interventions for tackling AI-generated disinformation.

First, adding more reliable, exogenous cues that can indicate the source and its quality may help people to assess the accuracy of information. Examples include adding flags for trustworthy sources (for example, a badge for verified users) or adding labels to AI-generated content that can warn or inform users. However, this raises several challenges, including which users to verify, how to appropriately label content without surfacing false-positives and how to detect that content is AI-generated in the first place. An alternative approach would be to mark the source of content upon creation. For example, some users who create AI-generated content might voluntarily flag their own output as AI-generated to avoid future decontextualization (deliberate or otherwise). Malicious actors, however, are unlikely to follow suit.

Second, leveraging collective intelligence can be effective in addressing not only human-generated but also AI-generated disinformation. Platforms could display social cues, for instance, that provide feedback on how many people actually shared a text or what amount of time people spent reading a text compared to the overall viewership. Such social cues are generally easy to implement and are already available for many platforms. However, they too have limitations, as AI-generated audio, images and videos are produced with the intent to capture attention, and adversarial actors often incorporate inauthentic engagement with their content to manipulate this very method of gauging legitimacy. Alternatively, crowdsourcing can be used by social media companies to assess content that is potentially misleading. For example, 'Community Notes' at X (formerly Twitter) enable users to add notes to content with contextualized information such as fact-checks. The success of such crowdsourced labels depends on the efforts of users, which may be low and thus require greater public buy-in or a platform-led incentive structure.

Third, boosting competencies can be effective in helping users to evaluate or verify information without relying on a single source.

Potentially valuable examples include media literacy, psychological inoculation, critical ignoring and lateral reading<sup>6</sup>. In particular, media literacy training will need to be adapted to the AI era. For example, one could train people in the basics of how generative AI models work or teach distinctive characteristics in their output, so that they better understand the potential risk of AI-generated disinformation. Yet, it is unclear whether this would increase general scepticism or actually improve truth discernment, and whether the AI-specific characteristics will remain consistent across models and as generative AI technology improves.

In sum, the combination of exogenous cues (which contribute to an online environment that cannot be gamed as easily) and the development of competencies that rely on multiple sources can provide a toolbox of behavioural interventions. Still, there are practical limitations behind the above behavioural interventions as the implementation generally depends on the willingness and capacity of social media platforms, which vary widely. Likewise, research is needed to understand how effective behavioural interventions will be in the era of AI-generated disinformation, and in comparison to existing strategies for increasing resilience to disinformation more generally.

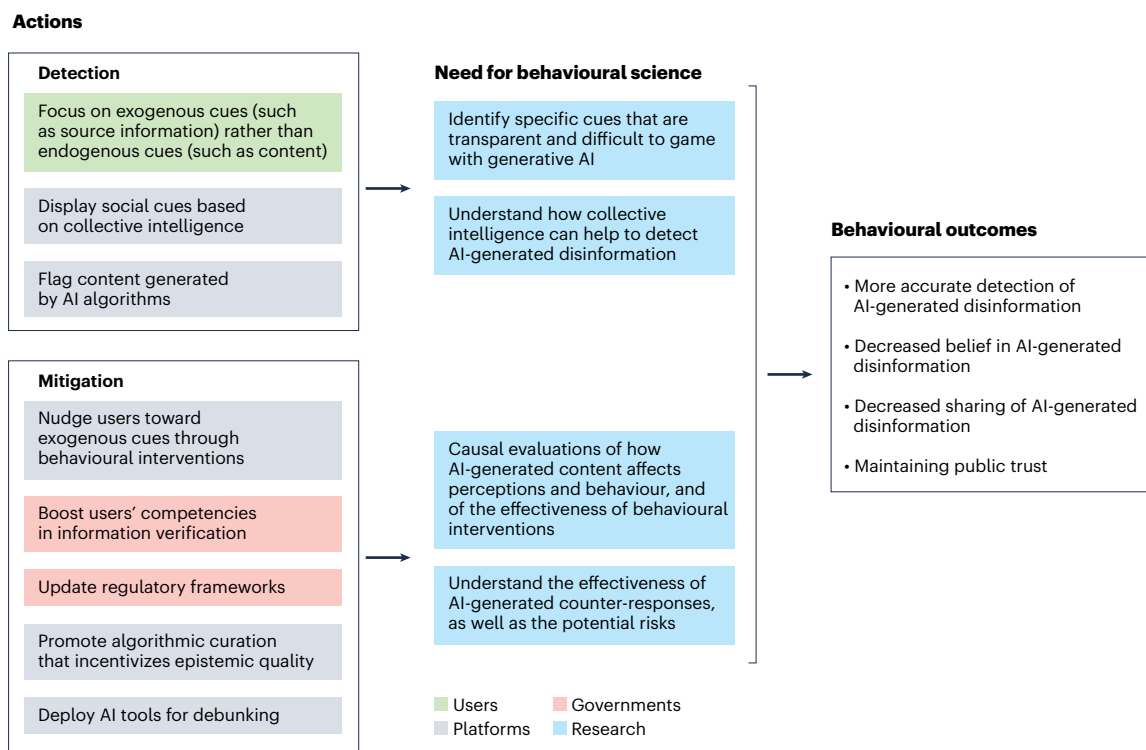
## Regulatory and technological mitigations

On the regulatory side, effective online governance is required that balances the importance of free speech against harms that may arise from AI-generated disinformation. Countries have begun developing relevant legal frameworks that prohibit or restrict the use of AI for generating content that is deceptive or manipulative<sup>7</sup>, although it is not clear how effective such regulation will be. Likewise, there are discussions that seek to limit the size of generative AI models<sup>8</sup>. However, such regulations may be sidestepped by malicious actors and current generative AI models may already be sufficient for producing high-quality content.

AI tools have also been developed that can be used to identify and flag AI-generated content in general or AI-generated disinformation in particular. However, these tools have notable shortcomings, such as surfacing false-positives. As generative AI models evolve, they will probably outpace detection tools to create a continuous cat-and-mouse game between generators and detectors. This highlights the need for continuous research and updates to keep detection tools effective and accurate – and to be aware when they are no longer reliable.

Technical solutions involve, for example, usage restrictions or fact-sensitive AI models as complements to regulation<sup>9</sup>. Further, watermarking could ease detection by embedding hidden signals in content to indicate that the content was produced by AI; these machine-readable signals enable social media platforms, which often serve as distribution sites, to recognize the content as AI-generated. Companies that produce AI tools for image and video generation could incorporate watermarks voluntarily as a form of self-regulation, even if not required to by law. However, there are enormous challenges behind coordinating watermarks across multiple digital platforms and across a growing number of AI tools (for example, through standardization).

Another technical solution is algorithmic amplification. Although AI-generated disinformation may be more persuasive, its effects are partially determined based on its reach. Depending on the extent to which algorithmic curation focuses on endogenous and exogenous cues (particularly as watermarking and provenance efforts progress), algorithmic curation may amplify or deamplify exposure to AI-generated disinformation. To this end, algorithmic curation that rewards cues of epistemic quality is desirable. For example, if the



**Fig. 1 | Role of behavioural science in the era of AI-generated disinformation.** Behavioural science can promote detection and mitigation strategies for humans in tackling AI-generated disinformation.

diversity of the readership is considered as a cue for algorithmic curation<sup>5</sup>, it could make it hard for disinformation to flourish – whether it is generated by AI or not.

Finally, research could enable AI technology itself to help to write counter-responses to disinformation (for example, a correction to fake news)<sup>10</sup>, which could be shared on social media (for example, via bots or other automated tools) and may help to prevent users from falling for or sharing it. A potential benefit of such counter-responses is that they are scalable, work in real-time and can be deployed by both social media platforms and third parties (for example, non-profit organizations or journalists) that seek to promote the accuracy of content posted on digital platforms. However, more research from behavioural science is needed to understand the effectiveness and ethics of AI-generated counter-responses, as well as the potential risks.

## Call for research

Generative AI technology has many positive applications, yet a negative externality is the democratization of disinformation content production: increasing its volume, velocity and potential persuasiveness while decreasing its cost. As we argue above, AI-generated disinformation challenges existing detection and mitigation strategies used by platforms and humans alike, and we therefore call for more research to update detection and mitigation strategies in light of AI-generated disinformation (Fig. 1). Impactful research questions seek to identify cues that are transparent and difficult to game with generative AI, to understand the effectiveness of behavioural interventions aimed at mitigating AI-generated disinformation, and to prepare new media

literacy training that is tailored to the upcoming challenges of the generative AI era. Eventually, to inform policy responses, we need rigorous, causal evaluations of how AI-generated content affects perceptions and behaviour, and of how detection and mitigation strategies could help.

**Stefan Feuerriegel**<sup>1,2</sup>✉, **Renée DiResta**<sup>3</sup>, **Josh A. Goldstein**<sup>4</sup>, **Srijan Kumar**<sup>5</sup>, **Philipp Lorenz-Spreen**<sup>6</sup>, **Michael Tomz**<sup>7,8</sup> & **Nicolas Pröllochs**<sup>9</sup>

<sup>1</sup>LMU Munich School of Management, LMU Munich, Munich, Germany.

<sup>2</sup>Munich Center for Machine Learning (MCML), Munich, Germany.

<sup>3</sup>Stanford Internet Observatory, Stanford University, Stanford, CA, USA.

<sup>4</sup>Center for Security and Emerging Technology, Georgetown University, Washington, DC, USA.

<sup>5</sup>College of Computing at Georgia Institute of Technology, Atlanta, GA, USA.

<sup>6</sup>Center for Adaptive Rationality, Max Planck Institute for Human Development, Berlin, Germany.

<sup>7</sup>Department of Political Science, Stanford University, Stanford, CA, USA.

<sup>8</sup>Stanford Institute for Economic Policy Research, Stanford University, Stanford, CA, USA.

<sup>9</sup>Department of Business and Economics, JLU Giessen, Giessen, Germany.

✉e-mail: [feuerriegel@lmu.de](mailto:feuerriegel@lmu.de)

Published online: 20 November 2023

## References

- Kreps, S. & Kriner, D. L. *New Media Soc.* <https://doi.org/10.1177/14614448231160526> (2023).
- Spitale, G., Biller-Andorno, N. & Germani, F. *Sci. Adv.* **9**, eadh1850 (2023).
- Kreps, S., McCain, R. M. & Brundage, M. *J. Exp. Political Sci.* **9**, 104–117 (2022).

- 
4. Lorenz-Spreen, P., Lewandowsky, S., Sunstein, C. R. & Hertwig, R. *Nat. Hum. Behav.* **4**, 1102–1109 (2020).
  5. Bhadani, S. et al. *Nat. Hum. Behav.* **6**, 495–505 (2022).
  6. Kozyreva, A., Lewandowsky, S. & Hertwig, R. *Psychol. Sci. Public Interest* **21**, 103–156 (2020).
  7. Hine, E. & Floridi, L. *Nat. Mach. Intell.* **4**, 608–610 (2022).
  8. Bender, E. M., Gebru, T., McMillan-Major, A. & Shmitchell, S. On the dangers of stochastic parrots: can language models be too big? In *FAccT '21: ACM Conf. on Fairness, Accountability, and Transparency*, 610–623 (ACM, 2021).
  9. Goldstein, J. A. et al. Preprint at *arXiv*, <https://doi.org/10.48550/arXiv.2301.04246> (2023).
  10. He, B., Ahamad, M. & Kumar, S. Reinforcement learning-based counter-misinformation response generation: a case study of COVID-19 vaccine misinformation. In *WWW '23: Proc. ACM Web Conf. 2023*, 2698–2709 (ACM, 2023).

## Competing interests

The authors declare no competing interests.

## Additional information

**Peer review information** *Nature Human Behaviour* thanks Jennifer Stromer-Galley and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.