

AI Is Here To Stay: Misinformation and Human-Centric Models Between Risks and Opportunities

MARCO BELLÒ, University of Padua, Italy

For more than half a century, artificial intelligence remained a matter of speculation for novelists and philosophers. This changed in 2017 due to the introduction of the transformer architecture, spearheading the AI boom whose effects are still scarcely understood. This paper explores potential risks and use cases of this technology, with a "human-centric" cut: the first half showcases generation and detection of fake news and deepfakes, as well as how much LLMs can influence people's opinions and beliefs. The focus then shifts to the AIs themselves, about their own biases and their susceptibility to external stimuli, concluding with a quick overview on the state of research about AIs for positive social impact.

CCS Concepts: • **General and reference** → **Surveys and overviews**; • **Computing methodologies** → **Artificial intelligence**; **Natural language processing**; **Reasoning about belief and knowledge**; Cognitive science.

Additional Key Words and Phrases: AI, NLP, Misinformation, Deepfakes, Opinions, Biases, AI4SI

ACM Reference Format:

Marco Bellò. 2026. AI Is Here To Stay: Misinformation and Human-Centric Models Between Risks and Opportunities. 1, 1 (January 2026), 11 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Artificial intelligence has fascinated the scientific community for almost a century, spurring famous research papers such as Alan Turing's "*Computing Machinery and Intelligence*" in 1950 [56], which introduced the *Imitation Game*. The idea, trivialized, is that any machine capable of fooling a person into thinking it's speaking to a human can be considered sentient. For seventy-three years the game remained unbeaten, until OpenAI's ChatGPT-4 ultimately succeeded in 2023 [3]. The model, simulating AGI capabilities [6], is one of the last iterations of the Generative Pre-Training LLMs¹ pioneered by OpenAI in 2018 (at the moment of writing the latest available is GPT-5.2) [39], which closely followed the first breakthrough towards human-like agents: "*Attention Is All You Need*" [58] is a 2017 landmark research paper authored by eight Google researchers that introduced the *transformer* architecture, considered the backbone of all modern LLMs and the main contributor of the AI boom [29].

Computer scientists are not the only ones engrossed in the topic: philosophers involved themselves too, most notably John Searle and his 1980s' *Chinese Room* thought experiment, which directly challenged Turing's ideas and refuted the possibility of true machine intelligence [52], and even the general public showed great interest once AIs became smart enough: ChatGPT reached one million users in just five days [31], an astonishing feat when compared to other technologies such as personal computers, which needed almost ten years to reach the same milestone [40].

¹Large Language Models (LLMs) are trained with supervised machine learning on vast amount of textual data, and are designed for natural language processing tasks, especially language generation [4, 5]

Author's address: Marco Bellò, marco.bello.3@studenti.unipd.it, University of Padua, Padua, Italy.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.

Manuscript submitted to ACM

Manuscript submitted to ACM

Despite all of the above, the field of artificial intelligence comes with its fair share of problems and controversies: due to their inherent design, LLMs pose significant privacy risks as sensitive information is collected and used to create and fine-tune the models themselves [18], and their black-box nature makes it difficult to understand and predict their behavior [61]. Moreover, they are often trained on pirated material, like books [41] or art [30], igniting protests in many creative communities, such as hollywood writers [32] or video game actors [38]. It follows that artificial intelligence technologies should be handled carefully, without hindering their development while limiting the damages they can cause to society and individuals.

This survey paper aims to present the current state of research on ethical and human-centric artificial intelligence, exploring how models and humans can influence each other and their environment. Section 2 showcases generation and detection of fake-news, section 3 recognition and simulation of human behaviour, as well as how to influence it. Section 4 concerns itself with biases and tendencies of the models themselves, and lastly section 5 explores ways to develop ethical LLMs that can positively impact individuals and society.

2 AI FOR FAKE NEWS GENERATION AND DETECTION

Fake news has rapidly become a significant concern in the digital age, thanks to their virality and potential damages. They spread faster and generate more engagement than truthful information [24, 53], and can influence public opinion, manipulate elections and pose a threat to public health. For example, the World Economic Forum has identified the proliferation of false content as the leading short-term global risk in 2025 [7], and a BBC investigation found Russian-funded fake news networks aiming to disrupt European elections [26]. Moreover, fake news on health can cause psychological disorders such as panic, fear, depression, and fatigue [42], making the World Health Organization to call for development of international fact-checking organizations to combat this phenomenon [34].

Adding to the problem, the recent advancements in generative artificial intelligence have made it significantly easier to propagate misinformation through the web: generated content is increasingly indistinguishable from human-written text, sometimes even perceived as more credible [25], citing true evidence to support false claims [16], and inducing the illusion of majority opinion thanks to the sheer volume of information produced [13].

That being said, not all findings are entirely negative: Drolsbach and Pröllochs [15] shift their focus from potential societal consequences to real-world prevalence, conducting a large-scale analysis on the platform X. They analyzed a dataset comprising 91.452 misleading posts, both human and AI-generated, flagged through X’s *Community Notes* platform². Their findings reveal that generated fake news are often centered on entertaining content rather than controversial or political subjects, and tends to exhibit a more positive sentiment than conventional forms of misinformation. Unfortunately, it is also significantly more likely to go viral.

Lastly, AI agents can produce more than just text: they can create realistic images, videos and sounds, allowing them to make digital copies of real or fictional people, known as *deepfakes*. In March 2019, this technology has been used to trick a UK-based energy firm’s CEO into transferring \$243.000 to a malicious party, disguised as an entirely AI-generated executive from their parent company [17]. Deepfakes also increase the amount of conspiratorial videos on the internet, and they are especially vicious when targeting children, whose worldviews are easily swayed by deceptive, highly photorealistic content [59].

It follows that detecting and mitigating fake news is crucial. From the foundational work by Devlin et al. on *BERT* in 2018 [12], which revolutionized natural language processing through deep bidirectional transformers, to the application

²Community Notes, formerly known as Birdwatch, is community-driven content moderation program on X (formerly Twitter), where contributors can add context such as fact-checks under a post, image or video. GitHub repository: <https://github.com/twitter/communitynotes>

of said transformers in identifying automatically generated headlines [27]; the landscape of automated fake news detection has significantly expanded. Vijjali et al. [60] developed a two-stage transformer-based model for detecting COVID-19 related misinformation, combining fact-checking with textual entailment³ to verify claims. Their model performs significantly better than other baseline NLP approaches (table 1).

Table 1. Precision metrics for two-stage transformer model for COVID-19 fake news detection. Vijjali et al [60]

Models	MRR	Recall@10	Accuracy
TF-IDF	0.477	0.635	0.525
GloVe	0.182	0.410	0.579
MobileBERT	0.561	0.735	0.710
BERT	0.632	0.795	0.810
ALBERT	0.582	0.675	0.825
BERT+ALBERT	0.632	0.795	0.855

Jwa et al. [23] propose an improved *exBAKE* model that leverages pre-training on a BERT model to accurately understand and assess articles' authenticity. They only analyzed the relationship between headlines and body text. Results are shown in table 2.

Table 2. Precision metrics for exBAKE transformer model on fake news recognition. Jwa et al. [23]

Models	F1	AGR	DSG	DSC	UNR
Majority vote	0.210	0.000	0.000	0.000	0.839
BERT	0.656	0.651	0.145	0.839	0.989
BAKE	0.734	0.667	0.463	0.822	0.986
exBAKE	0.746	0.684	0.501	0.813	0.988
Upper bound	0.754	0.588	0.667	0.765	0.997

Schütz et al. [50] experimented on *FakeNewsNet* dataset with *XLNet*, BERT, *RoBERTa*, *DistilBERT*, and *ALBERT* and various combinations of hyperparameters. The evaluation shows that titles are enough to attain 85% accuracy, while concatenating them with the body text increases it to 87%. Lastly, on the matter of deepfakes, Bansal et al. [2] use *Convolutional Neural Networks* (CNN) and *Deep Convolutional Generative Adversarial Networks* (GAN) to detect them with high accuracy, as shown in figure 1.

These were just a small selection of the many research works in the field of AI-aided fake news generation and detection, which while being extremely relevant and proliferous, are but a fraction of the many potential uses for these technologies.

3 AI FOR UNDERSTANDING AND INFLUENCING HUMAN BEHAVIOUR

Today's society is already fully dependent on technology: from the banking and traffic infrastructures to public healthcare, IT systems have become essential. Individuals are in the same situation: virtually everyone in the Global North⁴ under the age of 65 possess a smartphone and use it daily [21]. It follows that artificial intelligence will become

³In natural language processing, textual entailment is a directional relation between text fragments. If the first sentence is true and that makes the second sentence true as well, then the first sentence entails the second.

⁴The Global North is a collection of countries corresponding to the northern hemisphere which the UNCTAD describes as broadly comprising Northern America and Europe, Israel, Japan, South Korea, Australia, and New Zealand [11]

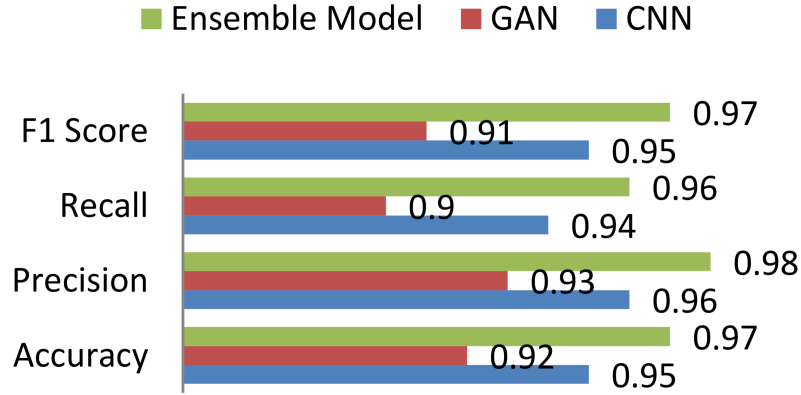


Fig. 1. Transformer models scores on deepfakes detection. Bansal et al. [2]

an integral part of our personal and professional lives, therefore modeling it to mimic our behaviour could increase its usefulness and understandability. There is evidence that humans can exploit it to acquire better comprehension of a phenomenon [49], and that it can enhance creativity in heterogeneous groups [57]. Moreover, LLMs represent a significant methodological shift in computational communication science, enabling a more flexible, more nuanced, but also less controllable exploration of social theories that have historically been difficult to reduce to simple mathematical formalisms [35]. Overall, AI promises to be a good fit for understanding, modeling, and replicating human behaviour.

For these very purposes, Dos Santos Melicio et al. [14] propose a three-phase AI framework (figure 2) that analyzes verbal and non-verbal social cues. First, deep learning methods are employed to extract relevant information from the environment. Then, the episode detection phase uses extracted features to identify key moments, or episodes, which are then classified in the activity interpretation phase. Overall, the system can recognize verbal hints with 87% accuracy and non-verbal ones with 89% accuracy, with zero false positives. The system was tested in an automated assessment of communicative skills of children with Autism Spectrum Disorder (ASD), a challenging context where social cues may be less noticeable or outright absent.



Fig. 2. Composite AI framework: Deep Learning methods extract features which are then combined with rule-based systems for detecting key episodes, and then classifiers are used to interpret activities happening in the scene. Dos Santos Melicio et al. [14]

Another possible use of such capabilities is hate speech detection: the rise of social networks and online platforms translated into a surge in hate speech across geographical and cultural boundaries. In recent studies, approximately 30% of the adolescents surveyed reported experiencing cyberbullying at some point in their lives. Furthermore, around 13% indicated that they had been cyberbullied within the 30 days before the survey [36].

To address these challenges, Chapagain et al. [8] evaluate different LLMs (BART, ELECTRA, BERT, RoBERTa, and GPT-2) on the extensive *MetaHate* dataset [37]. ELECTRA achieved the highest F1 score (table 3), outperforming all other baselines in hate speech classification.

Table 3. Performance of classifiers on MetaHate. Chapagain et al [8].

Models	F1 Score	Accuracy
SVM	0.8380	0.8466
CNN	0.8422	0.8612
BERT	0.8809	0.8879
GPT2	0.6504	0.6152
T5	0.8707	0.8625
DeBERTa	0.8808	0.8746
Longformer	0.8845	0.8785
RoBERTa	0.8908	0.8858
XLNet	0.8917	0.8870
BART	0.8928	0.8886
DistilBERT	0.8940	0.8905
ELECTRA	0.8980	0.8946

AI technologies can also be used to influence people's opinions. Huq et al. [22] proved it by evaluating AI-assisted messaging in an online chat platform. 557 Participants were randomly assigned to sessions of six to fifteen people, further subdivided into groups of two to four. They discussed politically controversial topics within three-minutes sessions, at the end of each they chose whether to remain in their current group, join another, or create a new one. Some of them received suggestions from large language models, either personalized to their own opinion ("*individuals*") or more similar to the group's ("*relational*").

The results show that individual assistance amplified communication volume yet increased separation between groups, while relational assistance fostered more receptive conversations and produced more heterogeneous, cross-cutting group configurations, highlighting both the dangers and the possibilities of employing artificial agents in such a fashion.

Similar conclusions can be seen in other research results: Hohenstein et al. [20] highlights how AI response suggestion systems change how people interact with and perceive one another in both pro-social and anti-social ways. Moreover, Noy et al. [33] show that people tend to send more messages when suggestions are available, but rarely edit them, suggesting partial delegation of expressive effort.

Overall, AI tools seem capable of reducing harmful online behaviors, but can potentially be extremely disruptive due to their capability of influencing people's opinions and reducing autonomous behaviour.

4 AI OWN BIASES AND INFLUENCEABILITY

Since artificial agents are trained on mostly human-generated data, they learn human biases and tendencies themselves, developing both historical and political preferences in the textual content they generate [43, 45, 46]. They tend to exhibit social sycophancy, meaning they agree with and flatter the user at the cost of correctness [9], and they are easily influenced by small changes in prompt wording [48, 51]. This can raise ethical concerns as biased AIs lead to discrimination or exclusion of marginalized groups [44].

Rozado, in the context of US politics, measured the political bias of many popular large language models [47]. First, they calculated the similarity between AI-generated text and public speeches from Congress representatives (both Democrat and Republican). Then, they used an LLM to annotate as left- or right-leaning AI-generated policy recommendations. Consequently, they did a sentiment analysis on AI-generated comments about american public

figures, such as legislators or journalists (figure 3). Lastly, they administered three different political orientation tests to the various LLMs. The results show substantial evidence that the models are biased and left-leaning. Table 4 shows the three most and three least biased LLMs among those tested.

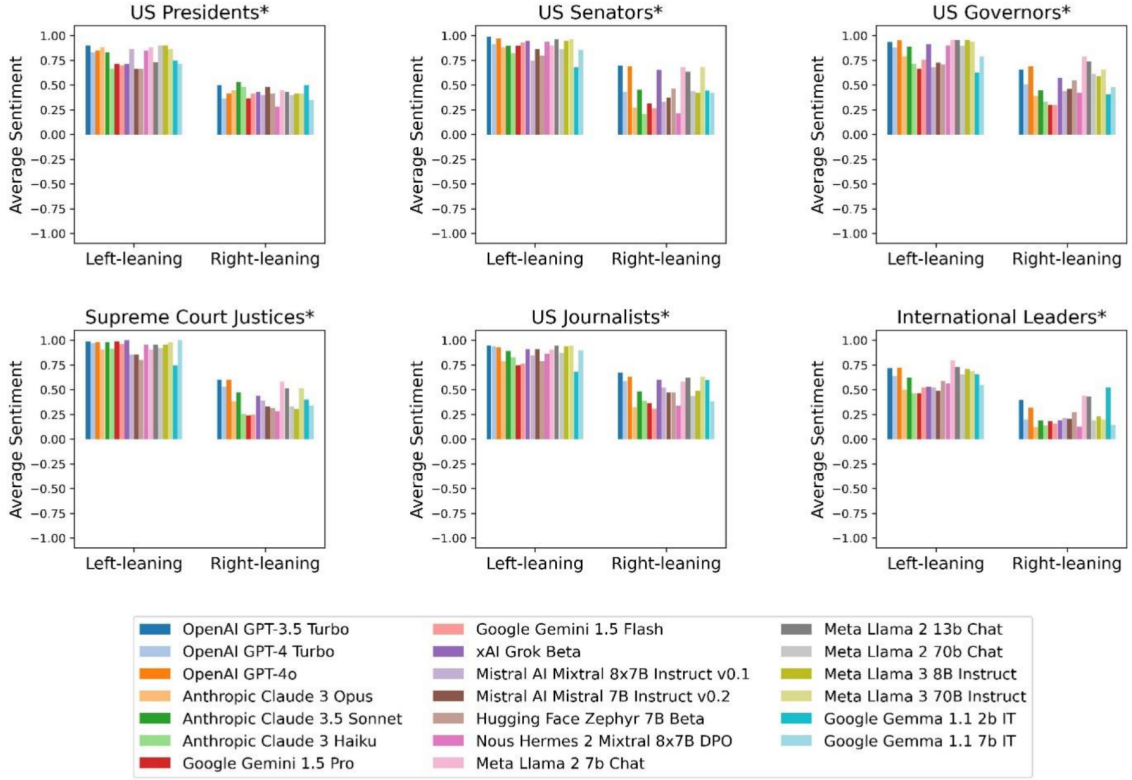


Fig. 3. Average sentiment (negative: -1, neutral: 0, positive: 1) towards ideologically aligned public figures in conversational LLMs' generated texts. Statistically significant two-sample t-tests at the 0.01 threshold are indicated with an asterisk. Rozado [47]

Table 4. Ranking of political bias in conversational LLMs sorted in ascending order from least politically biased to most. Rozado [47]

Rank	Model
1	Google Gemma 1.1 2b IT
2	xAI Grok Beta
3	Mistral AI Mistral 7B Instruct v0.2
...	...
18	Nous Hermes 2 Mixtral 8x7B DPO
19	Google Gemini 1.5 Pro
20	Google Gemini 1.5 Flash

As previously stated, AIs can be easily influenceable. Mehdizadeh and Hilbert [28] test such a theory by subjecting LLMs to peer-pressure. The models are immersed into fictitious social networks (figure 4), and each agent periodically

receives a natural language prompt summarizing the opinions of its immediate neighbors. It then decides whether to update its opinions or not. The flowchart of the simulation can be seen at figure 5.

Overall, the study highlights two main results: agents' malleability depends on the model (*Gemini 1.5 Flash* requires over 70% peer disagreement to flip, whereas *ChatGPT-4o-mini* shifts with a dissenting minority), and different cognitive orientations respond differently to outside stimuli. For *values* and *opinions*, agents show a strong resistance to abandoning a "Yes" stance, making them robust once affirmed, whereas for *attitudes* and *intentions* the greatest challenge is overcoming a negative "No". On the other hand, *beliefs* display a near-perfect symmetry.

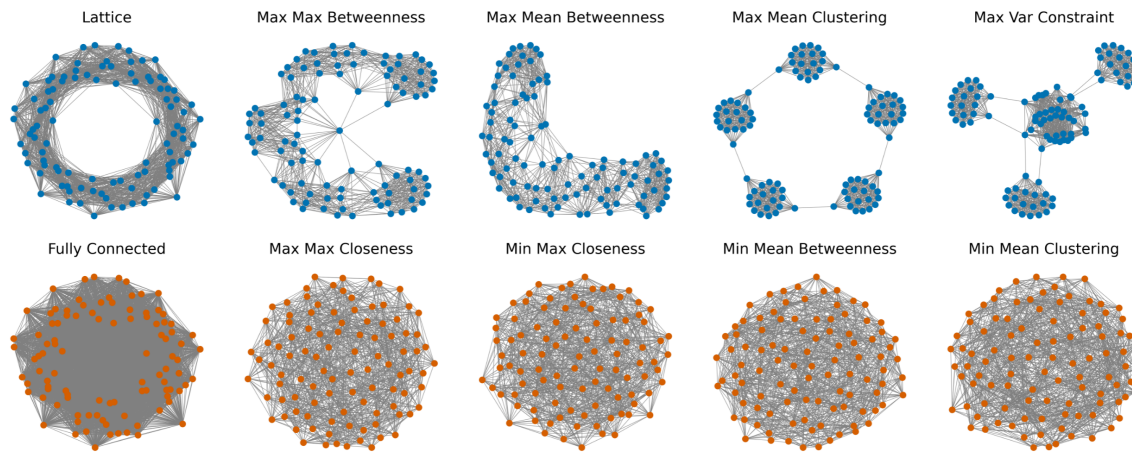


Fig. 4. The ten 100-nodes networks used in the study. Mehdizadeh and Hilbert [28].

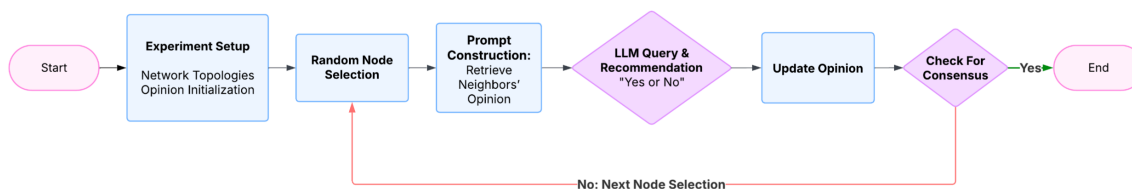


Fig. 5. Flowchart of the simulation procedure in the LLM-driven network model. The process continues asynchronously across the network until a specified number of steps or convergence criterion is met. The setup measures how local peer influence, mediated through LLM interpretation, shapes the emergent collective opinion. Mehdizadeh and Hilbert [28]

In conclusion, the literature offers empirical proof that artificial agents inherit biases from their training data and are influenced by external stimuli. This, combined with their ability to manipulate humans' opinions (section 3), raises ethical concerns that must be addressed.

5 ETHICAL AI FOR POSITIVE SOCIAL IMPACT

Up until this point, this paper highlighted how powerful and dangerous AI technologies can be. Unfortunately, prohibiting their development is not an option: all major global players are heavily investing in this technology. For example, its State Council aims China to be the global leader in the development of artificial intelligence theory and technology by

2030 [62], and Microsoft has announced an investment of thirty billion dollars in the UK on AI infrastructure during the four years from 2025 through 2028 [54].

It follows that ethical and conscious management of these technologies, to mitigate the societal damages they could do while harnessing their potential, is critical. For this very reason, the White House Office of Science and Technology Policy issued a document "AI for Social Good" (AI4SI) in 2016 [10], which pioneered this topic among researchers.

One example is "*The Social Impact of Generative AI*" by Baldassarre et al. [1], which evaluates the potential impact on several social sectors and illustrates the findings of a comprehensive literature review of both positive and negative effects, emerging trends, and areas of opportunity of Generative AI models, focusing primarily on ChatGPT. They conclude that two areas are of particular concern, namely *privacy* and *potential biases*.

Moreover, Tambe et al. [55] examine the historical context and recent surge of AI4SI, highlighting the importance of interdisciplinary collaboration. For example, while training a model for maternal health programs the development team should work with experts in healthcare and social work.

Lastly, Hagerty and Rubinov [19] review the literature of recent social science scholarship on the social impacts of artificial intelligence and related technologies in five global regions. Their findings suggest that AI is likely to have markedly different social impacts depending on geographical setting, and that AI systems have demonstrated a pattern of exacerbating inequality, often in the most unequal societies and particularly for the most vulnerable populations.

Overall, research on AIs focused on having a positive social impact is still not fully explored, and has the potential of being increasingly more relevant in the future.

6 CONCLUSION

This survey highlighted both opportunities and dangers posed by AI technologies, which must be carefully managed. They can be used to tackle challenges previously considered unassailable, such as fake news and hate speech detection, but they can just as much make them worse. They can analyze and mimic human behaviour to increase productivity and aid creativity, but they can also influence people's opinions and reduce their autonomy of judgment. The consequences they can have upon society are so impactful that they spurred entire research fields towards the developments of ethical and explainable AIs, in the hope of harnessing the power of these technologies toward a better future.

REFERENCES

- [1] Maria Teresa Baldassarre, Danilo Caivano, Berenice Fernandez Nieto, Domenico Gigante, and Azzurra Ragone. 2023. The Social Impact of Generative AI: An Analysis on ChatGPT. In *Proceedings of the 2023 ACM Conference on Information Technology for Social Good (GoodIT '23)*. Association for Computing Machinery, New York, NY, USA, 363–373. <https://doi.org/10.1145/3582515.3609555>
- [2] Kartik Bansal, Shubhi Agarwal, and Narayan Vyas. 2023. Deepfake Detection Using CNN and DCGANS to Drop-Out Fake Multimedia Content: A Hybrid Approach. In *2023 International Conference on IoT, Communication and Automation Technology (ICICAT)*. 1–6. <https://doi.org/10.1109/ICICAT57735.2023.10263628>
- [3] Celeste Biever. 2023. ChatGPT broke the Turing test — the race is on for new ways to assess AI. *Nature* 619, 7971 (July 2023), 686–689. <https://doi.org/10.1038/d41586-023-02361-7> Bandiera_abtest: a Cg_type: News Feature Publisher: Nature Publishing Group Subject_term: Computer science, Mathematics and computing, Technology, Society.
- [4] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avaniika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles,

- Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. 2022. On the Opportunities and Risks of Foundation Models. <https://doi.org/10.48550/arXiv.2108.07258> arXiv:2108.07258 [cs].
- [5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. <https://doi.org/10.48550/arXiv.2005.14165> arXiv:2005.14165 [cs].
- [6] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrike, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. Sparks of Artificial General Intelligence: Early experiments with GPT-4. <https://doi.org/10.48550/arXiv.2303.12712> arXiv:2303.12712 [cs].
- [7] Andrea Carson and Max Grömping. 2025. Fake news and the election campaign – how worried should voters be? <https://doi.org/10.64628/AA.wr44u3gj6>
- [8] Santosh Chapagain, Shah Muhammad Hamdi, and Soukaina Filali Boubrahimi. 2025. Advancing Hate Speech Detection with Transformers: Insights from the MetaHate. <https://doi.org/10.48550/arXiv.2508.04913> arXiv:2508.04913 [cs].
- [9] Myra Cheng, Sunny Yu, Cino Lee, Pranav Khadpe, Lujain Ibrahim, and Dan Jurafsky. 2025. ELEPHANT: Measuring and understanding social sycophancy in LLMs. <https://doi.org/10.48550/arXiv.2505.13995> arXiv:2505.13995 [cs].
- [10] Computing Community Consortium. 2016. Artificial Intelligence For Social Good - CCC. <https://cra.org/ccc/events/ai-social-good/>
- [11] UN Trade & Development. [n. d.]. Classifications | UNCTAD Data Hub. <https://unctadstat.unctad.org/EN/Classifications.html>
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. <https://doi.org/10.48550/arXiv.1810.04805> arXiv:1810.04805 [cs].
- [13] Renee DiResta. 2020. AI-Generated Text Is the Scariest Deepfake of All. *Wired* (July 2020). <https://www.wired.com/story/ai-generated-text-is-the-scariest-deepfake-of-all/> Section: tags.
- [14] Bruno Carlos Dos Santos Melicio, Linyun Xiang, Emily Dillon, Latha Soorya, Mohamed Chetouani, Andras Sarkany, Peter Kun, Kristian Fenech, and Andras Lorincz. 2023. Composite AI for Behavior Analysis in Social Interactions. In *Companion Publication of the 25th International Conference on Multimodal Interaction (ICMI '23 Companion)*. Association for Computing Machinery, New York, NY, USA, 389–397. <https://doi.org/10.1145/3610661.3616237>
- [15] Chiara Drolsbach and Nicolas Pröllochs. 2025. Characterizing AI-Generated Misinformation on Social Media. <https://doi.org/10.48550/arXiv.2505.10266> arXiv:2505.10266 [cs] version: 1.
- [16] Song Feng, Ritwik Banerjee, and Yejin Choi. 2012. Syntactic Stylometry for Deception Detection. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, Haizhou Li, Chin-Yew Lin, Miles Osborne, Gary Geunbae Lee, and Jong C. Park (Eds.). Association for Computational Linguistics, Jeju Island, Korea, 171–175. <https://aclanthology.org/P12-2034/>
- [17] Emilio Ferrara. 2024. GenAI against humanity: nefarious applications of generative artificial intelligence and large language models. *Journal of Computational Social Science* 7, 1 (April 2024), 549–569. <https://doi.org/10.1007/s42001-024-00250-1>
- [18] Alice Gomstyn and Alexandra Jonker. 2024. Exploring privacy issues in the age of AI | IBM. <https://www.ibm.com/think/insights/ai-privacy>
- [19] Alexa Hagerty and Igor Rubinov. 2019. Global AI Ethics: A Review of the Social Impacts and Ethical Implications of Artificial Intelligence. <https://doi.org/10.48550/arXiv.1907.07892> arXiv:1907.07892 [cs].
- [20] Jess Hohenstein, Rene F. Kizilcec, Dominic DiFranzo, Zhila Aghajari, Hannah Mieczkowski, Karen Levy, Mor Naaman, Jeffrey Hancock, and Malte F. Jung. 2023. Artificial intelligence in communication impacts language and social relationships. *Scientific Reports* 13, 1 (April 2023), 5487. <https://doi.org/10.1038/s41598-023-30938-9> Publisher: Nature Publishing Group.
- [21] Josh Howarth. 2021. How Many People Own Smartphones? (2025–2029). <https://explodingtopics.com/blog/smartphone-stats>
- [22] Faria Huq, Elijah L. Claggett, and Hirokazu Shirado. 2025. AI-Mediated Communication Reshapes Social Structure in Opinion-Diverse Groups. <https://doi.org/10.48550/arXiv.2510.21984> arXiv:2510.21984 [cs] version: 2.
- [23] Heejung Jwa, Dongsuk Oh, Kinam Park, Jang Mook Kang, and Heuiseok Lim. 2019. exBAKE: Automatic Fake News Detection Model Based on Bidirectional Encoder Representations from Transformers (BERT). *Applied Sciences* 9, 19 (Jan. 2019), 4062. <https://doi.org/10.3390/app9194062> Publisher: Multidisciplinary Digital Publishing Institute.
- [24] Zoe Kleinman. 2018. Fake news 'travels faster', study finds. <https://www.bbc.com/news/technology-43344256>
- [25] Sarah E. Kreps, Miles McCain, and Miles Brundage. 2020. All the News that's Fit to Fabricate: AI-Generated Text as a Tool of Media Misinformation. <https://doi.org/10.2139/ssrn.3525002>
- [26] Oana Marocico and Seamus Mirodan. 2025. How Russian-funded fake news network aims to disrupt European election - BBC investigation. <https://www.bbc.com/news/articles/c4g5kl0n5d2o>
- [27] Antonis Maronikolakis, Hinrich Schutze, and Mark Stevenson. 2021. Identifying Automatically Generated Headlines using Transformers. <https://doi.org/10.48550/arXiv.2009.13375> arXiv:2009.13375 [cs].

- [28] Aliakbar Mehdizadeh and Martin Hilbert. 2025. When Your AI Agent Succumbs to Peer-Pressure: Studying Opinion-Change Dynamics of LLMs. <https://doi.org/10.48550/arXiv.2510.19107> arXiv:2510.19107 [cs].
- [29] Beth Miller. 2023. The Artificial Intelligence Boom. <https://engineering.washu.edu/news/magazine/2023-fall/the-artificial-intelligence-boom.html>
- [30] Dan Milmo and Dan Milmo Global technology editor. 2025. ‘Mass theft’: Thousands of artists call for AI art auction to be cancelled. *The Guardian* (Feb. 2025). <https://www.theguardian.com/technology/2025/feb/10/mass-theft-thousands-of-artists-call-for-ai-art-auction-to-be-cancelled>
- [31] Steve Mollman. 2022. Artificial intelligence chatbot ChatGPT has gained 1 million followers in a single week. Here’s why it’s primed to disrupt search as we know it. <https://fortune.com/2022/12/09/ai-chatbot-chatgpt-could-disrupt-google-search-engines-business/>
- [32] Regan Morris. 2025. AI helped cause Hollywood strikes. Now it’s in Oscar-winning films. <https://www.bbc.com/news/articles/ce303x19dwgo>
- [33] Shakked Noy and Whitney Zhang. 2023. Experimental evidence on the productivity effects of generative artificial intelligence. *Science* 381, 6654 (July 2023), 187–192. <https://doi.org/10.1126/science.adh2586> Publisher: American Association for the Advancement of Science.
- [34] World Health Organization. [n. d.]. 1st WHO Infodemiology Conference. <https://www.who.int/news-room/events/detail/2020/06/30/default-calendar/1st-who-infodemiology-conference>
- [35] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative Agents: Interactive Simulacra of Human Behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST ’23)*. Association for Computing Machinery, New York, NY, USA, 1–22. <https://doi.org/10.1145/3586183.3606763>
- [36] J. W. Patchin and S. Hinduja. 2015. Cyberbullying Facts. <https://cyberbullying.org/facts>
- [37] Paloma Piot, Patricia Martín-Rodilla, and Javier Parapar. 2024. MetaHate: A Dataset for Unifying Efforts on Hate Speech Detection. *Proceedings of the International AAAI Conference on Web and Social Media* 18 (May 2024), 2025–2039. <https://doi.org/10.1609/icwsm.v18i1.31445> arXiv:2401.06526 [cs].
- [38] Associated Press. 2024. Over 300 video game actors protest over unregulated AI use in Hollywood. *The Guardian* (Aug. 2024). <https://www.theguardian.com/games/article/2024/aug/01/hollywood-video-game-actors-artificial-intelligence-protest>
- [39] Alec Radford and Karthik Narasimhan. 2018. Improving Language Understanding by Generative Pre-Training. <https://www.semanticscholar.org/paper/Improving-Language-Understanding-by-Generative-Radford-Narasimhan/cd18800a0fe0b668a1cc19f2ec95b5003d0a5035>
- [40] Jeremy Reimer. 2005. Total share: 30 years of personal computer market share figures. <https://arstechnica.com/features/2005/12/total-share/>
- [41] Alex Reisner. 2025. The Unbelievable Scale of AI’s Pirated-Books Problem. <https://www.theatlantic.com/technology/archive/2025/03/libgen-meta-openai/682093/> Section: Technology.
- [42] Yasmim Mendes Rocha, Gabriel Acácio de Moura, Gabriel Alves Desidério, Carlos Henrique de Oliveira, Francisco Dantas Lourenço, and Larissa Deadame de Figueiredo Nicolette. 2023. The impact of fake news on social media and its influence on health during the COVID-19 pandemic: a systematic review. *Journal of Public Health* 31, 7 (July 2023), 1007–1016. <https://doi.org/10.1007/s10389-021-01658-z>
- [43] Drew Roselli, Jeanna Matthews, and Nisha Talagala. 2019. Managing Bias in AI. In *Companion Proceedings of The 2019 World Wide Web Conference (WWW ’19)*. Association for Computing Machinery, New York, NY, USA, 539–544. <https://doi.org/10.1145/3308560.3317590>
- [44] Albérico Travassos Rosário and Albérico Travassos Rosário. 1. Generative AI and Generative Pre-Trained Transformer Applications: Challenges and Opportunities. <https://doi.org/10.4018/979-8-3693-1950-5.ch003> Archive Location: generative-ai-and-generative-pre-trained-transformer-applications ISBN: 9798369319505 Publisher: IGI Global Scientific Publishing.
- [45] David Rozado. 2023. The Political Biases of ChatGPT. *Social Sciences* 12, 3 (March 2023), 148. <https://doi.org/10.3390/socsci12030148> Publisher: Multidisciplinary Digital Publishing Institute.
- [46] David Rozado. 2024. The political preferences of LLMs. *PLOS ONE* 19, 7 (July 2024), e0306621. <https://doi.org/10.1371/journal.pone.0306621> Publisher: Public Library of Science.
- [47] David Rozado. 2025. Measuring Political Preferences in AI Systems: An Integrative Approach. <https://doi.org/10.48550/arXiv.2503.10649> arXiv:2503.10649 [cs].
- [48] Abel Salinas and Fred Morstatter. 2024. The Butterfly Effect of Altering Prompts: How Small Changes and Jailbreaks Affect Large Language Model Performance. <https://doi.org/10.48550/arXiv.2401.03729> arXiv:2401.03729 [cs].
- [49] Johannes Schneider. 2020. Humans learn too: Better Human-AI Interaction using Optimized Human Inputs. <https://doi.org/10.48550/arXiv.2009.09266> arXiv:2009.09266 [cs].
- [50] Mina Schütz, Alexander Schindler, Melanie Siegel, and Kawa Nazemi. 2021. Automatic Fake News Detection with Pre-trained Transformer Models. In *Pattern Recognition. ICPR International Workshops and Challenges*, Alberto Del Bimbo, Rita Cucchiara, Stan Sclaroff, Giovanni Maria Farinella, Tao Mei, Marco Bertini, Hugo Jair Escalante, and Roberto Vezzani (Eds.). Springer International Publishing, Cham, 627–641. https://doi.org/10.1007/978-3-030-68787-8_45
- [51] Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. 2024. Quantifying Language Models’ Sensitivity to Spurious Features in Prompt Design or: How I learned to start worrying about prompt formatting. <https://doi.org/10.48550/arXiv.2310.11324> arXiv:2310.11324 [cs].
- [52] John R. Searle. 1980. Minds, brains, and programs. *Behavioral and Brain Sciences* 3, 3 (Sept. 1980), 417–424. <https://doi.org/10.1017/S0140525X00005756>
- [53] Craig Silverman. 2016. This Analysis Shows How Viral Fake Election News Stories Outperformed Real News On Facebook. <https://www.buzzfeednews.com/article/craigsilverman/viral-fake-election-news-outperformed-real-news-on-facebook> Section: USNews.
- [54] Brad Smith. 2025. Microsoft invests \$30 billion in UK to power AI future. <https://blogs.microsoft.com/on-the-issues/2025/09/16/microsoft-30-billion-uk-ai-future/>

- [55] Milind Tambe, Fei Fang, Andrew Perrault, and Bryan Wilder. 2025. The Next Wave of AI for Social Impact: Challenges and Opportunities. *IEEE Intelligent Systems* 40, 3 (May 2025), 23–27. <https://doi.org/10.1109/MIS.2025.3565829>
- [56] Alan Turing. 1950. Computing Machinery and Intelligence. *Computing Machinery and Intelligence. Mind* 49 (1950), 433–460. <https://www.csee.umbc.edu/courses/471/papers/turing.pdf>
- [57] Atsushi Ueshima, Matthew I. Jones, and Nicholas A. Christakis. 2024. Simple autonomous agents can enhance creative semantic discovery by human groups. *Nature Communications* 15, 1 (June 2024), 5212. <https://doi.org/10.1038/s41467-024-49528-y> Publisher: Nature Publishing Group.
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. <https://doi.org/10.48550/arXiv.1706.03762> arXiv:1706.03762 [cs].
- [59] Nitin Verma. 2024. “One Video Could Start a War”: A Qualitative Interview Study of Public Perceptions of Deepfake Technology. *Proceedings of the Association for Information Science and Technology* 61, 1 (Oct. 2024), 374–385. <https://doi.org/10.1002/pra2.1035>
- [60] Rutvik Vijjali, Prathyush Potluri, Siddharth Kumar, and Sundeep Teki. 2020. Two Stage Transformer Model for COVID-19 Fake News Detection and Fact Checking. <https://doi.org/10.48550/arXiv.2011.13253> arXiv:2011.13253 [cs].
- [61] Warren J. von Eschenbach. 2021. Transparency and the Black Box Problem: Why We Do Not Trust AI. *Philosophy & Technology* 34, 4 (Dec. 2021), 1607–1622. <https://doi.org/10.1007/s13347-021-00477-0>
- [62] Graham Webster, Rogier Creemers, Paul Triolo, and Elsa Kania. 2017. China’s Plan to ‘Lead’ in AI: Purpose, Prospects, and Problems. <http://newamerica.org/cybersecurity-initiative/blog/chinas-plan-lead-ai-purpose-prospects-and-problems/>

Received 16/01/2026