

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

3 MARCO BELLÒ, University of Padua, Italy

4 **ACM Reference Format:**

5 Marco Bellò. 2025. AI Is Here To Stay: Misinformation and Human-Centric Models Between Risks and Opportunities. 1, 1 (January 2025),
6 4 pages. <https://doi.org/XXXXXX.XXXXXXXX>

7 **1 INTRODUCTION**

8 Artificial intelligence has fascinated the scientific community for almost a century, spurring famous research papers
9 such as Alan Turing's "*Computing Machinery and Intelligence*" in 1950 [27], which introduced the *imitation game*. The
10 idea, trivialized, is that any machine capable of fooling a person into thinking it's speaking to a human can be considered
11 sentient. For seventy-three years the game remained unbeaten, until OpenAI's ChatGPT-4 ultimately succeeded in
12 2023 [2]. The model, simulating AGI capabilities [5], is one of the last iterations of the Generative Pre-Training LLMs¹
13 pioneered by OpenAI in 2018 (at the moment of writing the latest available is GPT-5.2) [21], which closely followed
14 the first breakthrough towards human-like agents: "*Attention Is All You Need*" [28] is a 2017 landmark research paper
15 authored by eight Google researchers that introduced the *transformer* architecture, considered the backbone of all
16 modern LLMs and the main contributor of the AI boom [15].

17 Computer scientists are not the only ones engrossed in the topic: philosophers involved themselves too, most notably
18 Jhon Searle and his 1980s' *chinese room* thought experiment, which directly challenged Turing's ideas and refuted
19 the possibility of true machine intelligence [25], and even the general public showed great interest once AIs became
20 smart enough: ChatGPT reached one million users in just five days [17], an astonishing feat when compared to other
21 technologies such as personal computers, which needed almost ten years to reach the same milestone [22].

22 Despite all of the above, the field of artificial intelligence comes with its fair share of problems and controversies:
23 due to their inherent design, LLMs pose significant privacy risks as sensitive information is collected and used to create
24 and fine-tune the models themselves [11], and their black-box nature makes it difficult to understand and predict their
25 behavior [30]. Moreover, they are often trained on pirated material, like books [23] or art [16], igniting protests in many
26 creative communities, such as hollywood writers [18] or video game actors [20]. It follows that artificial intelligence
27 technologies should be handled carefully, without hindering their development while limiting the damages they can
28 cause to society and individuals.

29
30
31
32
33
34
35
36
37
38
39
40
41
42
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 [3, 4]

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

44
45 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
46 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
47 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
48 servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

49 © 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
50 Manuscript submitted to ACM

51
52 Manuscript submitted to ACM

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 have rapidly become a significant concern in the modern digital age, thanks to their virality and potential damages. They spread faster and generate more engagement than truthful information [12, 26], and can influence public opinion, manipulate elections and pose a threat to public health: the European Union issued guidelines to online platforms and search engines to mitigate the impact on misinformation on elections [1], the World Economic forum has identified the proliferation of false content as the leading short-term global risk in 2025 [6], and a BBC investigation found Russian-funded fake news networks aiming to disrupt european elections [14]. Moreover, fake news on health can cause psychological disorders and panic, fear, depression, and fatigue [24], and the World Health Organization called for the development of international fact-checking organizations to combat this phenomenon [19].

Adding to the problem, the recent advancements in generative artificial intelligence have made it significantly easier to propagate disinformation throughout the web: generated content is increasingly indistinguishable from human-written text, sometimes even perceived as more credible [13], citing true evidence to support false claims [9], and inducing the illusion of majority opinion thanks to the sheer volume of information produced [8]. Some works highlights how

Artificial agents can generate more than just text: they can create realistic images, videos and sounds, allowing them to reproduce digital twins of real or fictional people, known as deepfakes. In March 2019, such a technology has been used to trick a UK-based energy firm’s CEO into transferring \$243.000 to a convincingly mimicked company’s German parent firm’s CEO [10]. Deepfakes also increased the amount of conspiratorial videos on the internet, and they are especially vicious when targeting children, whose worldviews are easily swayed by deceptive—and highly photorealistic—content [29].

It follows that detecting and mitigating fake news is crucial, especially since the rise of AI-generated content has made disinformation easier to spread and more convincing. From the foundational work by Devlin et al. on *BERT* in 2018 [7], which revolutionized natural language processing trough deep bidirectional transformers, to innovative detection models like *exBAKE* and the application of transformers [TODO]

3 AI ON HUMANS

4 AI OWN BIASES

5 ETHICAL AI

REFERENCES

- [1] Federico Baccini. [n. d.]. *Against fake news and misinformation in European elections*. <https://www.eunews.it/en/2024/03/26/eu-commission-steps-up-work-against-online-fake-news-ahead-of-european-elections/> Section: Net & Tech.
- [2] Celeste Biiever. 2023. ChatGPT broke the Turing test — the race is on for new ways to assess AI. 619, 7971 (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.
- [3] 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,

- 105 Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh,
 106 Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter
 107 Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth
 108 Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal
 109 Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir
 110 Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles,
 111 Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher
 112 Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori
 113 Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang,
 114 William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi
 115 Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. [n. d.]. On the Opportunities and Risks of Foundation Models.
<https://doi.org/10.48550/arXiv.2108.07258> arXiv:2108.07258 [cs]
- 116 [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish
 117 Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
 118 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark,
 119 Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. [n. d.]. Language Models are Few-Shot Learners. <https://doi.org/10.48550/arXiv.2005.14165> arXiv:2005.14165 [cs]
- 120 [5] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott
 121 Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. [n. d.]. Sparks of Artificial General Intelligence: Early experiments with
 122 GPT-4. <https://doi.org/10.48550/arXiv.2303.12712> arXiv:2303.12712 [cs]
- 123 [6] Andrea Carson and Max Grömping. [n. d.]. *Fake news and the election campaign – how worried should voters be?* <https://doi.org/10.64628/AA.wr44u3gi6>
- 124 [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. [n. d.]. BERT: Pre-training of Deep Bidirectional Transformers for Language
 125 Understanding. <https://doi.org/10.48550/arXiv.1810.04805> arXiv:1810.04805 [cs]
- 126 [8] Renee DiResta. [n. d.]. AI-Generated Text Is the Scariest Deepfake of All. ([n. d.]). <https://www.wired.com/story/ai-generated-text-is-the-scariest-deepfake-of-all/> Section: tags.
- 127 [9] Song Feng, Ritwick Banerjee, and Yejin Choi. [n. d.]. Syntactic Stylometry for Deception Detection. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)* (Jeju Island, Korea, 2012-07), Haizhou Li, Chin-Yew Lin, Miles Osborne, Gary Geunbae Lee, and Jong C. Park (Eds.). Association for Computational Linguistics, 171–175. <https://aclanthology.org/P12-2034/>
- 128 [10] Emilio Ferrara. [n. d.]. GenAI against humanity: nefarious applications of generative artificial intelligence and large language models. 7, 1 ([n. d.]),
 129 549–569. <https://doi.org/10.1007/s42001-024-00250-1>
- 130 [11] Alice Gomstyn and Alexandra Jonker. [n. d.]. *Exploring privacy issues in the age of AI* / IBM. <https://www.ibm.com/think/insights/ai-privacy>
- 131 [12] Zoe Kleinman. [n. d.]. *Fake news ‘travels faster’*, study finds. <https://www.bbc.com/news/technology-43344256>
- 132 [13] Sarah E. Kreps, Miles McCain, and Miles Brundage. [n. d.]. All the News that’s Fit to Fabricate: AI-Generated Text as a Tool of Media Misinformation.
<https://doi.org/10.2139/ssrn.3525002>
- 133 [14] Oana Marocico and Seamus Mirodan. [n. d.]. *How Russian-funded fake news network aims to disrupt European election - BBC investigation.*
<https://www.bbc.com/news/articles/c4g5kl0n5d2o>
- 134 [15] Beth Miller. [n. d.]. *The Artificial Intelligence Boom.* <https://engineering.washu.edu/news/magazine/2023-fall/the-artificial-intelligence-boom.html>
- 135 [16] Dan Milmo and Dan Milmo Global technology editor. [n. d.]. ‘Mass theft’: Thousands of artists call for AI art auction to be cancelled. ([n. d.]).
<https://www.theguardian.com/technology/2025/feb/10/mass-theft-thousands-of-artists-call-for-ai-art-auction-to-be-cancelled>
- 136 [17] Steve Mollman. [n. d.]. *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/>
- 137 [18] Regan Morris. [n. d.]. *AI helped cause Hollywood strikes. Now it’s in Oscar-winning films.* <https://www.bbc.com/news/articles/ce303x19dwgo>
- 138 [19] 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>
- 139 [20] Associated Press. [n. d.]. Over 300 video game actors protest over unregulated AI use in Hollywood. ([n. d.]). <https://www.theguardian.com/games/article/2024/aug/01/hollywood-video-game-actors-artificial-intelligence-protest>
- 140 [21] Alec Radford and Karthik Narasimhan. [n. d.]. Improving Language Understanding by Generative Pre-Training. <https://www.semanticscholar.org/paper/Improving-Language-Understanding-by-Generative-Radford-Narasimhan/cd18800a0fe0b668a1cc19f2ec95b5003d0a5035>
- 141 [22] Jeremy Reimer. [n. d.]. *Total share: 30 years of personal computer market share figures.* <https://arstechnica.com/features/2005/12/total-share/>
- 142 [23] Alex Reisner. [n. d.]. *The Unbelievable Scale of AI’s Pirated-Books Problem.* <https://www.theatlantic.com/technology/archive/2025/03/libgen-meta-openai/682093/> Section: Technology.
- 143 [24] Yasmim Mendes Rocha, Gabriel Acácio de Moura, Gabriel Alves Desidério, Carlos Henrique de Oliveira, Francisco Dantas Lourenço, and
 144 Larissa Deadame de Figueiredo Nicolete. [n. d.]. The impact of fake news on social media and its influence on health during the COVID-19
 145 pandemic: a systematic review. 31, 7 ([n. d.]), 1007–1016. <https://doi.org/10.1007/s10389-021-01658-z>
- 146 [25] John R. Searle. [n. d.]. Minds, brains, and programs. 3, 3 ([n. d.]), 417–424. <https://doi.org/10.1017/S0140525X00005756>

- 157 [26] Craig Silverman. [n. d.]. *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.
- 158 [27] Alan Turing. [n. d.]. Computing Machinery and Intelligence. 49 ([n. d.]), 433–460. <https://www.csee.umbc.edu/courses/471/papers/turing.pdf>
- 159 [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention Is
160 All You Need. arXiv:1706.03762 [cs.CL] <https://arxiv.org/abs/1706.03762>
- 161 [29] Nitin Verma. [n. d.]. “One Video Could Start a War”: A Qualitative Interview Study of Public Perceptions of Deepfake Technology. 61, 1 ([n. d.]),
162 374–385. <https://doi.org/10.1002/pra2.1035>
- 163 [30] Warren J. von Eschenbach. [n. d.]. Transparency and the Black Box Problem: Why We Do Not Trust AI. 34, 4 ([n. d.]), 1607–1622. <https://doi.org/10.1007/s13347-021-00477-0>
- 164
- 165
- 166
- 167
- 168
- 169
- 170
- 171
- 172
- 173
- 174
- 175
- 176
- 177
- 178
- 179
- 180
- 181
- 182
- 183
- 184
- 185
- 186
- 187
- 188
- 189
- 190
- 191
- 192
- 193
- 194
- 195
- 196
- 197
- 198
- 199
- 200
- 201
- 202
- 203
- 204
- 205
- 206
- 207
- 208