

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

3
4 MARCO BELLÒ, University of Padua, Italy
5
6

7 **ACM Reference Format:**

8 Marco Bellò. 2025. AI Is Here To Stay: Misinformation and Human-Centric Models Between Risks and Opportunities. 1, 1 (January 2025),
9 3 pages. <https://doi.org/XXXXXX.XXXXXXXX>

10 **1 INTRODUCTION**

11 Artificial intelligence has fascinated the scientific community for almost a century, spurring famous research papers
12 such as Alan Turing's "*Computing Machinery and Intelligence*" in 1950 [15], which introduced the *imitation game*. The
13 idea, trivialized, is that any machine capable of fooling a person into thinking it's speaking to a human can be considered
14 sentient. For seventy-three years the game remained unbeaten, until OpenAI's ChatGPT-4 ultimately succeeded in
15 2023 [1]. The model, simulating AGI capabilities [4], is one of the last iterations of the Generative Pre-Training LLMs¹
16 pioneered by OpenAI in 2018 (at the moment of writing the latest available is GPT-5.2) [11], which closely followed
17 the first breakthrough towards human-like agents: "*Attention Is All You Need*" [16] is a 2017 landmark research paper
18 authored by eight Google researchers that introduced the *transformer* architecture, considered the backbone of all
19 modern LLMs and the main contributor of the AI boom [6].

20 Computer scientists are not the only ones engrossed in the topic: philosophers involved themselves too, most notably
21 Jhon Searle and his 1980s' *chinese room* thought experiment, which directly challenged Turing's ideas and refuted
22 the possibility of true machine intelligence [14], and even the general public showed great interest once AIs became
23 smart enough: ChatGPT reached one million users in just five days [8], an astonishing feat when compared to other
24 technologies such as personal computers, which needed almost ten years to reach the same milestone [12].

25 Despite all of the above, the field of artificial intelligence comes with its fair share of problems and controversies:
26 due to their inherent design, LLMs pose significant privacy risks as sensitive information is collected and used to create
27 and fine-tune the models themselves [5], and their black-box nature makes it difficult to understand and predict their
28 behavior [17]. Moreover, they are often trained on pirated material, like books [13] or art [7], igniting protests in many
29 creative communities, such as hollywood writers [9] or video game actors [10]. It follows that artificial intelligence
30 technologies should be handled carefully, without hindering their development while limiting the damages they can
31 cause to society and individuals.

32
33
34
35
36
37
38
39
40
41¹Large Language Models (LLMs) are trained with supervised machine learning on vast amount of textual data, and are designed for natural language
42 processing tasks, especially language generation [2, 3]

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

3 AI ON HUMANS

4 AI OWN BIASES

5 ETHICAL AI

REFERENCES

- [1] Celeste Biever. 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.
- [2] 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, Avanika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforshan, 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. [n. d.]. On the Opportunities and Risks of Foundation Models. <https://doi.org/10.48550/arXiv.2108.07258> arXiv:2108.07258 [cs]
- [3] 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. [n. d.]. Language Models are Few-Shot Learners. <https://doi.org/10.48550/arXiv.2005.14165> arXiv:2005.14165 [cs]
- [4] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. [n. d.]. Sparks of Artificial General Intelligence: Early experiments with GPT-4. <https://doi.org/10.48550/arXiv.2303.12712> arXiv:2303.12712 [cs]
- [5] Alice Gomstyn and Alexandra Jonker. [n. d.]. *Exploring privacy issues in the age of AI / IBM*. <https://www.ibm.com/think/insights/ai-privacy>
- [6] Beth Miller. [n. d.]. *The Artificial Intelligence Boom*. <https://engineering.washu.edu/news/magazine/2023-fall/the-artificial-intelligence-boom.html>
- [7] 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>
- [8] 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/>
- [9] Regan Morris. [n. d.]. *AI helped cause Hollywood strikes. Now it's in Oscar-winning films*. <https://www.bbc.com/news/articles/ce303x19dwgo>
- [10] 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>
- [11] 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>
- [12] Jeremy Reimer. [n. d.]. *Total share: 30 years of personal computer market share figures*. <https://arstechnica.com/features/2005/12/total-share/>
- [13] Alex Reisner. [n. d.]. *The Unbelievable Scale of AI’s Pirated-Books Problem*. <https://www.theatlantic.com/technology/archive/2025/03/libgen-metopenai/682093/> Section: Technology.

- 105 [14] John R. Searle. [n. d.]. Minds, brains, and programs. 3, 3 ([n. d.]), 417–424. <https://doi.org/10.1017/S0140525X00005756>
106 [15] Alan Turing. [n. d.]. Computing Machinery and Intelligence. 49 ([n. d.]), 433–460. <https://www.csee.umbc.edu/courses/471/papers/turing.pdf>
107 [16] Ashish Vaswani, Noah Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is
108 All You Need. arXiv:1706.03762 [cs.CL] <https://arxiv.org/abs/1706.03762>
109 [17] 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>

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156