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

MARCO BELLÒ, University of Padua, Italy

## 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 [32], 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 LLMs[1] pioneered by OpenAI in 2018 (at the moment of writing the latest available is GPT-5.2) [25], which closely followed the first breakthrough towards human-like agents: *"Attention Is All You Need"* [33] 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 [19].

Computer scientists are not the only ones engrossed in the topic: philosophers involved themselves too, most notably Jhon Searle and his 1980s' *chinese room* thought experiment, which directly challenged Turing's ideas and refuted the possibility of true machine intelligence [30], and even the general public showed great interest once AIs became smart enough: ChatGPT reached one million users in just five days [21], an astonishing feat when compared to other technologies such as personal computers, which needed almost ten years to reach the same milestone [26].

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 [13], and their black-box nature makes it difficult to understand and predict their behavior [36]. Moreover, they are often trained on pirated material, like books [27] or art [20], igniting protests in many creative communities, such as hollywood writers [22] or video game actors [24]. 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.

---

[1]Large 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.

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 digital age, thanks to their virality and potential damages. They spread faster and generate more engagement than truthful information [15, 31], 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 [7], and a BBC investigation found Russian-funded fake news networks aiming to disrupt european elections [17]. Moreover, fake news on health can cause psychological disorders and panic, fear, depression, and fatigue [28], and the World Health Organization called for the development of international fact-checking organizations to combat this phenomenon [23].

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 [16], citing true evidence to support false claims [11], and inducing the illusion of majority opinion thanks to the sheer volume of information produced [9]. That being said, not all findings are entirely negative: Drolsbach and Pröllochs [10] 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 trough X's *Community Notes* platform [2]. 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, such a 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 [12]. 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 [34].

It follows that detecting and mitigating fake news is crucial. From the foundational work by Devlin et al. on *BERT* in 2018 [8], which revolutionized natural language processing trough deep bidirectional transformers, to the application of said transformers in identifying automatically generated headlines, significantly outperforming humans, by Maronikolakis et al. [18], the landscape of automated fake news detection has significantly expanded. Vijjali et al. [35] developed a two-stage transformer-based model for detecting COVID-19 related misinformation, combining fact-checking with textual entailment to verify claims. Their model performs significantly better than other baseline NLP approaches (table 1). Jwa et al. [14] 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

---

[2]Community Notes, formerly known as Birdwatch, is community-drive 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

body. Results can be seen in table 2. Schütz et al. [29] experimented on *FakeNewsNet* dataset with *XLNet*, BERT, *RoBERTa*, *DistilBERT*, and *ALBERT* and various combinations of hyperparameters. The evaluation shows that already short texts are enough to attain 85% accuracy on the test set. Using the body text and a concatenation of both reach up to 87% accuracy. 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.
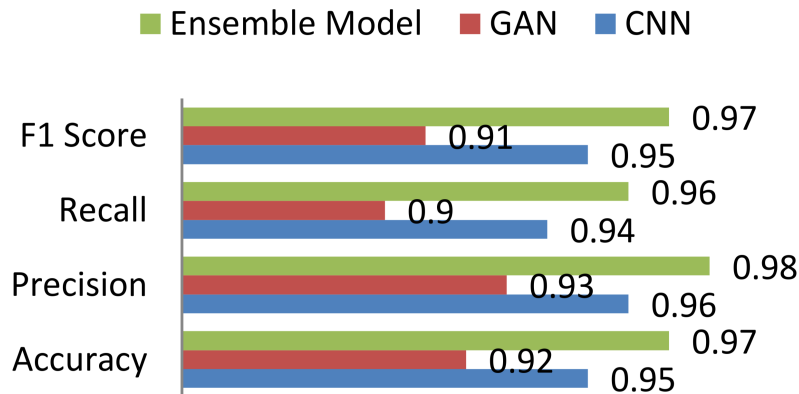


Fig. 1. Transformer models scores on deepfakes detection, Bansal et al. [2]

Table 1. Precision metrics for two-stage transformer model for COVID-19 fake news detection, Vijjali et al [35]

| 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 |

Table 2. Precision metrics for exBAKE transformer model on fake news recognition, Jwa et al. [14]

| 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 |

## 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] Kartik Bansal, Shubhi Agarwal, and Narayan Vyas. [n. d.]. 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)* (2023-06). 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. 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.

[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, Avanika 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. [n. d.]. 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. [n. d.]. 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 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]

[7] Andrea Carson and Max Grömping. [n. d.]. *Fake news and the election campaign – how worried should voters be?* https://doi.org/10.64628/AA.wr44u3gj6

[8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. [n. d.]. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. https://doi.org/10.48550/arXiv.1810.04805 arXiv:1810.04805 [cs]

[9] 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.

[10] Chiara Drolsbach and Nicolas Pröllochs. [n. d.]. Characterizing AI-Generated Misinformation on Social Media. https://doi.org/10.48550/arXiv.2505.10266 arXiv:2505.10266 [cs] version: 1.

[11] Song Feng, Ritwik 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/

[12] Emilio Ferrara. [n. d.]. GenAI against humanity: nefarious applications of generative artificial intelligence and large language models. 7, 1 ([n. d.]), 549–569. https://doi.org/10.1007/s42001-024-00250-1

[13] Alice Gomstyn and Alexandra Jonker. [n. d.]. *Exploring privacy issues in the age of AI | IBM.* https://www.ibm.com/think/insights/ai-privacy

[14] Heejung Jwa, Dongsuk Oh, Kinam Park, Jang Mook Kang, and Heuiseok Lim. [n. d.]. exBAKE: Automatic Fake News Detection Model Based on Bidirectional Encoder Representations from Transformers (BERT). 9, 19 ([n. d.]), 4062. https://doi.org/10.3390/app9194062 Publisher: Multidisciplinary Digital Publishing Institute.

[15] Zoe Kleinman. [n. d.]. *Fake news 'travels faster', study finds.* https://www.bbc.com/news/technology-43344256

[16] 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

[17] 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

[18] Antonis Maronikolakis, Hinrich Schutze, and Mark Stevenson. [n. d.]. Identifying Automatically Generated Headlines using Transformers. https://doi.org/10.48550/arXiv.2009.13375 arXiv:2009.13375 [cs]

[19] Beth Miller. [n. d.]. *The Artificial Intelligence Boom.* https://engineering.washu.edu/news/magazine/2023-fall/the-artificial-intelligence-boom.html

[20] 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

[21] 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/

[22] Regan Morris. [n. d.]. *AI helped cause Hollywood strikes. Now it's in Oscar-winning films.* https://www.bbc.com/news/articles/ce303x19dwgo

[23] 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

[24] Associated Press. [n. d.]. Over 300 video game actors protest over unregulated AI use in Hollywood. ([n. d.]). https://www.theguardian.acm/games/article/2024/aug/01/hollywood-video-game-actors-artificial-intelligence-protest

[25] 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

[26] Jeremy Reimer. [n. d.]. *Total share: 30 years of personal computer market share figures.* https://arstechnica.com/features/2005/12/total-share/

[27] 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.

[28] Yasmim Mendes Rocha, Gabriel Acácio de Moura, Gabriel Alves Desidério, Carlos Henrique de Oliveira, Francisco Dantas Lourenço, and Larissa Deadame de Figueiredo Nicolete. [n. d.]. The impact of fake news on social media and its influence on health during the COVID-19 pandemic: a systematic review. 31, 7 ([n. d.]), 1007–1016. https://doi.org/10.1007/s10389-021-01658-z

[29] Mina Schütz, Alexander Schindler, Melanie Siegel, and Kawa Nazemi. [n. d.]. Automatic Fake News Detection with Pre-trained Transformer Models. In *Pattern Recognition. ICPR International Workshops and Challenges* (Cham, 2021), Alberto Del Bimbo, Rita Cucchiara, Stan Sclaroff, Giovanni Maria Farinella, Tao Mei, Marco Bertini, Hugo Jair Escalante, and Roberto Vezzani (Eds.). Springer International Publishing, 627–641. https://doi.org/10.1007/978-3-030-68787-8_45

[30] John R. Searle. [n. d.]. Minds, brains, and programs. 3, 3 ([n. d.]), 417–424. https://doi.org/10.1017/S0140525X00005756

[31] 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.

[32] Alan Turing. [n. d.]. Computing Machinery and Intelligence. 49 ([n. d.]), 433–460. https://www.csee.umbc.edu/courses/471/papers/turing.pdf

[33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. arXiv:1706.03762 [cs.CL] https://arxiv.org/abs/1706.03762

[34] Nitin Verma. [n. d.]. "One Video Could Start a War": A Qualitative Interview Study of Public Perceptions of Deepfake Technology. 61, 1 ([n. d.]), 374–385. https://doi.org/10.1002/pra2.1035

[35] Rutvik Vijjali, Prathyush Potluri, Siddharth Kumar, and Sundeep Teki. [n. d.]. 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]

[36] 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