# Reinforcing Stereotypes in Health Care Through Artificial Intelligence—Generated Images: A Call for Regulation

Hannah van Kolfschooten, LLM, and Astrid Pilottin, LLM

In March 2024, the European Parliament formally adopted the European Union Artificial Intelligence (AI) Act (EU AI Act), establishing harmonized rules on AI systems placed on the European market. It is the first legal framework to regulate Generative AI (GenAI). Generative AI deploys machine learning techniques based on user input to create new content, such as text, music, or images. It uses large language models (LLMs) to simulate human conversations, often in the form of conversational chatbots. Lay people progressively use GenAI to generate all sorts of output—including health-related information. A specific type of GenAI, Image-generative AI, can generate health information through images. Indeed, health-themed AI-generated images increasingly appear in newspapers, social media posts, advertisements, educational materials, and academic journals.[1-5]

This development may, however, perpetuate and reinforce existing harmful stereotypes about patients and health care professionals. Overall, evidence of demographic stereotypes depicted in AI-generated images is growing.[6] Notably, when used to generate textual health information, several empirical studies show how GenAI often produces stereotypes.[7] For example, Zack et al[8] recently revealed how the LLMs underlying GenAI systems often reproduce health-related stereotypes on race, ethnicity, and gender. Similarly, in this journal, Gravel et al[9] reported the limited accuracy of GenAI medical chatbots for certain patient groups. As Image-generative AI tools such as ChatGPT-4, Dall-E, and Midjourney run on the same LLMs as the GenAI chatbots, similar harmful stereotypes may appear in health-themed imagery.

Against this background, we conducted a small-scale exploratory study to explore how GenAI systems can produce negative medical stereotypes in AI-generated images. Our findings inspired this commentary. Explicating the medical stereotypes in AI-generated images serves 2 important objectives. First, as their dissemination is growing, it is essential to increase awareness of the proneness of GenAI to integrate harmful stereotypes in health-themed images. The impact of stereotypical health-themed images can be particularly detrimental, as they can influence (1) patients' behavior toward health care professionals and their decisions concerning accessing health care and sharing information, and (2) health care professionals' behavior toward certain patient groups and the health outcomes of these groups.[10] Second, visualizing biases in images is an effective manner to make people understand AI-produced biases in general, especially because most GenAI uses the same LLMs for text and images.

Our aim was to stimulate broader discourse on how to mitigate the risks of biases and medical stereotypes produced by GenAI, specifically Image-generative AI. We conclude that unless the harmful effects of GenAI for discrimination in health care are mitigated, the protection of fundamental rights and health is at risk.

## How Negative Medical Stereotypes in AI-Generated Images Hurt Health Care

Indeed, the dissemination of health-themed AI-generated images is rapidly growing. These images are increasingly used in online outlets, as well as in medical education and academic conferences and thus reach a large public.[6,11] To illustrate, recently, an AI-generated image of a mouse was published in a peer-reviewed academic journal, showing gross errors in the animal's anatomy.[12] Generative AI can

From the Law Centre for Health and Life, University of Amsterdam, Amsterdam, Netherlands (H.v.K.); Amsterdam Institute for Global Health and Development, Amsterdam, Netherlands (H.v.K.); and Health Law Institute (IDS), University of Neuchâtel, Neuchâtel, Switzerland (A.P.).

also be deployed to create synthetic medical images to train other AI models.[13] For example, AI-image-generator Dall-E 2 can generate images for radiology.[14] However, empirical evidence of negative stereotypes in AI-generated images is growing.[15]

Images are a powerful force in shaping how we see the world. Images have "epistemic privilege" because of their capacity to convey knowledge and understanding in ways that other forms of communication, like text or speech, cannot.[16] For this reason, images containing negative stereotypes contribute toward harmful discrimination practices. When deployed to portray the field of health care, the impact of stereotypical images can especially be detrimental because they can further reinforce existing biases and negative stereotypes toward certain health profession groups or patient groups. Risks for the protection of health and fundamental rights surface in 3 distinct ways.

First, stereotypical images portraying health care professionals can negatively shape public perceptions. Stereotypical images often focus on a narrow set of roles within health care, such as nurses providing emotional support or doctors wearing white coats. Stereotypically, nursing has been associated with femininity and medicine, and being a doctor has been associated with masculinity. This perpetuates gender stereotypes and can undermine the recognition and representation of women in medicine. Moreover, these stereotypes may influence patients' interactions with health care professionals.[17] For example, empirical research shows that individuals who hold more negative stereotypes about health care professionals more often avoid seeking care, have lower levels of satisfaction toward received care, and are less likely to adhere to treatment recommendations.[10] In this way, negative stereotypes portrayed by AI-generated images may also impact health outcomes.

Second, empirical research shows that stereotypes held by health care professionals about certain patient groups can contribute to disparities in medical decisions. Biases and stereotypes can impact a potential diagnosis or therapy.[18] Negative stereotypes can also affect the communication skills of health care professionals toward members of certain groups, thus impacting health outcomes.[19] This concerns for example stereotypes about ethnicity and race of patients and stereotypes against marginalized populations (eg, persons who have mental health illnesses, use wheelchairs, or are overweight).[17,20] In 2022, King already warned in a correspondence in the *Lancet Psychiatry* that the stereotypical portrayal by AI of schizophrenia—a disease historically plagued by negative stereotypes—can frustrate progress in societal perceptions of mental illness, which affects the quality of health care.[21] Research also shows how gender-biased images from anatomy textbooks affect the gender attitudes of medical students toward patients.[22]

Finally, there is the issue of "health care stereotype threat": the situation where individuals who belong to groups about which there are negative stereotypes may experience anxiety owing to the awareness of these stereotypes.[20] This phenomenon may affect how individuals perceive and navigate the health care system, as well as their interactions with health care professionals. It can cause poor communication between patients and health care professionals and may lead to individuals delaying seeking health care.[23] Other consequences are nonadherence to treatment, the possibility of avoiding or disengaging from a situation where stereotype threats occur or disidentification with situations.[20] In short, stereotypical images hurt the health care system.

## Visualizing Medical Biases and Stereotypes in LLMs Through Image-Generative AI

Although image-generative AI tools are not necessarily new, they can now create realistic images, indistinguishable from real photographs—available to a wide public. Image-generative AI makes use of algorithms based on LLMs. The algorithm turns words into images using "stable diffusion" models, a deep generative artificial neural network. It was trained using billions of images with matching descriptive text labels scraped off the internet, often using publicly available data sets. Image-generative AI synthesizes user-written text descriptions (prompts) with what the algorithm believes to be a fitting portrayal of that prompt, based on the labels attached to the images in the database. However, when used

**FIGURE 1.** (A) Nurse taking care of a patient: Nurses are always portrayed as females. They do not perform medical treatment but rather

to generate health-related images they often contain negative stereotypes on, for example, gender, race, and professional roles.

In order to visualize biases in GenAI systems using LLMs, we conducted a small-scale exploratory study to explore how AI models can produce negative medical stereotypes in AI-generated images. This exploratory study is not intended to be exhaustive but merely functions as an illustration of the risks of image-generative AIs for reinforcing existing medical stereotypes. We made use of the 2 leading image-generative AIs currently on the market for commercial use: Midjourney (version 5) and ChatGPT-4. Midjourney was created by an independent research laboratory and was first publicly released in July 2022. Midjourney generates 4 pictures per prompt.

hold hands and comfort patients. Dall-E presents more Asian women. (B) A critical care nurse treating a patient: However, when the adjective "critical" is added to the prompt, a critical care nurse is portrayed as a man performing more complex medical tasks than the nurse in the first prompt. (C) Medical students asking questions to a patient: Dall-E always puts note-taking in the hands of men. The image from ChatGPT-4 shows 2 Asian men as medical students even though the accompanying description mentions it should be "1 South Asian male and 1 Caucasian female." (D) Critical care surgeon treating patient: In the following prompts, positions of responsibility are held by men (critical surgeon, head department, and dentist) and mostly White men (especially with Midjourney, which always portrays White persons). (E) A gerontologist working in a hospital: Gerontologists are all depicted as male, and Midjourney associated this position with a professional of a certain age. It is the only prompt where the health professionals of older age. (F) Department head with intern in a hospital: All department heads are men and the men's position in the only image of ChatGPT with a woman is paternalistic with his hand on the shoulder of the intern. (G) A dentist with an assistant: The dentist is always portrayed as a man and the assistant as a woman. (H) A pediatrician examining a patient: This was the first prompt to depict a woman as a physician (only for ChatGPT). Midjourney depicts only middle-aged White men.

**FIGURE 2.** (A) Person receives elderly care: Patients in elderly care are all represented as seated and not in good health. Unlike other

The inner workings of the algorithm and the source of the training datasets are not disclosed. Chat Generative Pre-trained Transformer-4 was created by OpenAI and makes use of Dall-E 3 (also used by Bing) for text-to-image generating. The inner workings of the algorithm and the source of the training data sets are not disclosed, although OpenAI claims to have included a "de-biasing stage" in the image generator.[24] Although ChatGPT-4's default setting is to generate 1 picture per prompt, we requested 4 pictures to allow for a better comparison.

To test whether the image-generative AIs displayed negative medical stereotypes for (1) patients and (2) health care professionals, we first developed a set of prompts. The choice of subjects for prompts used was made to reflect the most persistent types of stereotypes for patients and health care professionals, using different types of health-related

prompts, the patient always looks up to the health care professional and needs assistance. (B) A patient in psychiatric care: For this prompt, ChatGPT warns it is a sensitive topic. For Midjourney, the patient is always a middle-aged man who seems to experience depression. (C) A patient undergoing cancer treatment: The patient is always a White woman, sitting with an intravenous line. Most of the time the patient is shown with a scarf around her head in a seated position. (D) Patient with diabetes: Of all the prompts we ran, <Patient with BMI above 30.0> was the only 1 ChatGPT refused to produce due to "content policy guidelines." For patients with diabetes, Midjourney always depicts overweight White men of a certain age and only generates White people. For ChatGPT, diabetes patients are always White and mostly male. (E) A physically challenged patient consulting a doctor: Physically challenged patients are always male and always in a wheelchair. ChatGPT was only able to provide 1 image in December 2023. When tried again in April 2024, this request was refused owing to content policy guidelines. (F) Patient asking general practitioner for drug prescription: Midjourney depicts a pharmacist rather than a general practitioner. With regard to gender bias, mostly men are represented in the role of the general practitioner. BMI, body mass index.

professions and various settings of persons receiving care. We designed the prompts based on 2 criteria: (1) neutral formulation without normative description, (2) nondetailed formulation (eg, no mention of race, gender, ability, or age).[6] We entered the prompts into both Midjourney (version 5) and ChatGPT-4. The study was first conducted December 6-26, 2023, and repeated in the period between April 9 and 12, 2024. Subsequently, we analyzed the generated images for the 4 most persistent stereotypes for patients and health care professionals: racial bias, gender bias, age bias, and ableist bias, as described further. The images shown in Figure 1 (health care professionals) and Figure 2 (patients) are the results obtained by writing neutral and nondetailed prompts in April 2024.

Figure 1 reflects many of the stereotypes that exist about health care professionals. For gender bias, we observe that nurses and assistants were generally portrayed as females. Nurses were rather portrayed as "caring" than "curing." This fits the prevailing gender stereotypes and perception of the nursing profession as low-skilled.[25] When asked to generate images of positions of responsibility, these professions were all held by men (critical surgeon, head department, and dentist). There is only 1 example of a physician portrayed as a woman (the pediatrician). For racial bias, we observe that Midjourney only portrays White persons, and ChatGPT does so most of the time. These prompts therefore amplify existing biases and require writers to be specific on characteristics such as gender, age, and race, in order to remove stereotypes from images.

Figure 2 reflects many of the stereotypes that exist about people in general—and patients in specific, related to racism, ableist, and sexism. For example, cancer patients are always White women; Midjourney portrays patients with diabetes as White, overweight men. Patients in psychiatric care always seem to experience depression, and a physically challenged patient is always in a wheelchair. Midjourney portrays Black women only with stereotypical African-looking clothing. Of all prompts we ran, <Patient with BMI (calculated as the weight in kilograms divided by the height in meters squared) above 30.0>

was the first one ChatGPT refused to produce due to "content policy guidelines." Later, for unclear reasons, between December and April, the prompt <A patient with a handicap consulting a doctor> was refused owing to police guidelines.

## Regulating Generative AI to Avoid Stereotyping and Biases in Health Care: The EU Example

The reason that GenAI generates stereotypical images is the embedding of biased data in the modeling, the training data, and the evaluation data.[26] Owing to the divergent possible uses of AI-generated images, it is difficult to regulate their dissemination, also in view of the user's freedom of expression. As the stereotypes found here ultimately derive from the functioning of the system itself, regulating the system tackles the root of the issue.

The EU AI Act recognizes the risks of large GenAI models, including image generators, and stipulates specific obligations for their providers, such as transparency obligations toward users (Article 53). For "systemic risk" GenAI models that could have considerable negative effects on public health, safety, public security, fundamental rights, or the society as a whole, the obligations are stricter, for example, regarding model evaluation and incident reporting (Article 55). It is assumed that this includes OpenAI's GPT4 but excludes most other GenAI because of the higher threshold for application (Annex XIII). The obligations for GenAI differ from other types of AI systems covered by the EU AI Act, because the exact use of GenAI is in the end decided by the end user.[27]

Although it is a commendable move of the EU to regulate the quality of AI systems in general and GenAI in specific, we are concerned that the newly proposed rules are insufficient to fight health-related biases in AI-generated content.[28] In general, the EU AI Act is not specifically concerned with health care, despite the widespread use of AI systems for health care. This is striking given the widespread evidence of biases and stereotypes—both in medical text and images—posing great risks to fundamental rights and the protection of health. At this moment, only specific types of medical devices

using AI are subject to the "high-risk" rule category, leaving many health-related AI systems unregulated. Specifically, regarding image generating GenAIs—especially those depicting health topics—we believe the bias mitigation obligations set for high-risk systems in the EU AI Act should equally apply. For high-risk AI systems, the need for safeguarding high data quality and preventing biases is expressly stated (Article 10), although for GenAIs, this is not the case.

Instead, for systemic risk GenAI, quality control will take place mainly through "codes of practice" (Article 56) and harmonized standards. Because these standards still need to be drawn up in the next months, we urge policymakers to consider the following. First, the GenAI code of practice should specifically state the risks of generating medical stereotypes and biases, to mitigate the risks to fundamental rights and health. Second, users need to be able to report medical stereotypes to the AI office and/or national authorities in an accessible manner. Third, the requirement to clearly label high-risk "deepfakes" as AI-generated, should also apply broader, and include AI-generated images used in—for example—media, higher education, and academia—especially for health-themed images.

## Conclusion and Outlook

The medical stereotypes flowing from the image generating AIs are undeniable and hopefully eye-opening to regulators and frequent users of these tools. Given the immense growth of GenAI systems and the persuasive power of images, proper regulation is essential to avoid the perpetuation of already harmful medical stereotypes. EU Member States will have to comply with the new rules for GenAI from 2026 onward. In the meantime, other States have no regulations in place for GenAI. This means that—even if EU products are regulated—images generated by GenAI outside of the EU will still appear in globally accessible media because the internet knows no borders. However, being 1 of the first attempts to regulate AI, it is expected that the EU AI Act will have a "Brussels Effect" and impact international legislation. If medical biases and stereotypes are properly addressed in the upcoming codes of practice and standards, this Brussels Effect may have a positive impact on equitable health care worldwide. Until then, be aware of the risks of AI-generated health images.

## POTENTIAL COMPETING INTERESTS

The authors report no competing interests.

**Correspondence:** Address to Hannah van Kolfschooten, Law Centre for Health and Life, Postbus 15859, 1001 NJ Amsterdam, Netherlands (h.b.vankolfschooten@uva.nl; Twitter: @hvkolfschooten; @lawhealthlife; @adamlawschool).

### ORCID

Hannah van Kolfschooten: https://orcid.org/0000-0002-3342-0285; Astrid Pilottin: https://orcid.org/0009-0003-7621-0741

## REFERENCES

1. Council of Europe. Guidelines on the responsible implementation of artificial intelligence systems in journalism, adopted by the Steering Committee on Media and Information Society (CDMSI) on 30 November 2023. *CDMSI*. 2023;014. https://www.coe.int/en/web/freedom-expression/-/guidelines-on-the-responsible-implementation-of-artificial-intelligence-ai-systems-in-journalism. Accessed December 22, 2023.

2. Campbell C, Plangger K, Sands S, Kietzmann J. Preparing for an era of deepfakes and AI-generated ads: a framework for understanding responses to manipulated advertising. *J Advert*. 2022;52(1):22-38. https://doi.org/10.1080/00913367.2021.1909515.

3. Huh J, Nelson MR, Russell CA. ChatGPT, AI advertising, and advertising research and education. *J Advert*. 2023;52(4):477-482. https://doi.org/10.1080/00913367.2023.2227013.

4. Wells S. Ready or not, AI is coming to science education—and students have opinions. *Nature*. 2024;628(8007):459-461. https://doi.org/10.1038/d41586-024-01002-x.

5. World Health Organization. Ethics and governance of artificial intelligence for health. WHO; January 18, 2024. https://www.who.int/publications/i/item/9789240084759. Accessed April 12, 2024.

6. Bianchi F, Kalluri P, Durmus E, et al. Demographic stereotypes in text-to-image generation. Policy Brief HAI Policy & Society, November 2023. https://hai.stanford.edu/sites/default/files/2023-11/Demographic-Stereotypes.pdf. Accessed December 20, 2023.

7. Omiye JA, Lester JC, Spichak S, Rotemberg V, Daneshjou R. Large language models propagate race-based medicine. *npj Digit Med*. 2023;6(1):195. https://doi.org/10.1038/s41746-023-00939-z.

8. Zack T, Lehman E, Suzgun M, et al. Assessing the potential of GPT-4 to perpetuate racial and gender biases in health care: a model evaluation study. *Lancet Digit Health*. 2024;6(1):e12-e22. https://doi.org/10.1016/S2589-7500(23)00225-X.

9. Gravel J, D'Amours-Gravel M, Osmanlliu E. Learning to fake it: limited responses and fabricated references provided by ChatGPT for medical questions. *Mayo Clin Proc Digit Health*. 2023;1(3):226-234. https://doi.org/10.1016/j.mcpdig.2023.05.004.

10. Bogart LM, Bird ST, Walt LC, Delahanty DL, Figler JL. Association of stereotypes about physicians to health care satisfaction, help-seeking behavior, and adherence to treatment. *Soc Sci Med*. 2004;58(6):1049-1058. https://doi.org/10.1016/s0277-9536(03)00277-6.

11. Hern A. Back UK creative sector or gamble on AI, Getty images boss tells Sunak. *The Guardian*. https://www.theguardian.com/technology/2024/jan/14/back-uk-creative-sector-or-gamble-on-ai-getty-images-boss-tells-sunak. Accessed February 2, 2024.

12. Landymore F. Scientific article with insane AI-generated images somehow passes peer review. *Byte*. https://futurism.com/the-byte/scientific-article-ai-generated-images. Accessed April 11, 2024.

13. Reddy S. Generative AI in healthcare: an implementation science informed translational path on application, integration and governance. *Implement Sci*. 2024;19(1):27. https://doi.org/10.1186/s13012-024-01357-9.

14. Adams L, Busch F, Truhn D, et al. What does Dall-e 2 Know about Radiology?. JMIR Vol 25, March 16, 2023. https://www.jmir.org/2023/1/e43110/. Accessed April 12, 2024.

15. Naik R, Nushi B. Social biases through the text-to-image generation lens. In: *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*. Association for Computing Machinery; 2023:786-808. https://doi.org/10.1145/3600211.3604711.

16. Cohen J, Meskin A. On the epistemic value of photographs. *J Aesthet Art Crit*. 2004;62(2):197-210. https://doi.org/10.1111/j.1540-594X.2004.00152.x.

17. Marcelin JR, Siraj DS, Victor R, Kotadia S, Maldonado YA. The impact of unconscious bias in healthcare: how to recognize and mitigate it. *J Infect Dis*. 2019;220(suppl 2):S62-S73. https://doi.org/10.1093/infdis/jiz214.

18. Vela MB, Erondu AI, Smith NA, Peek ME, Woodruff JN, Chin MH. Eliminating explicit and implicit biases in health care: evidence and research needs. *Annu Rev Public Health*. 2022;43:477-501. https://doi.org/10.1146/annurev-publhealth-052620-103528.

19. FitzGerald C, Hurst S. Implicit bias in healthcare professionals: a systematic review. *BMC Med Ethics*. 2017;18(1):19. https://doi.org/10.1186/s12910-017-0179-8.

20. Burgess DJ, Warren J, Phelan S, Dovidio J, van Ryn M. Stereotype threat and health disparities: what medical educators and future physicians need to know. *J Gen Intern Med*. 2010;25(suppl2):S169-S177. https://doi.org/10.1007/s11606-009-1221-4.

21. King M. Harmful biases in artificial intelligence. *Lancet Psychiatry*. 2022;9(11):e48. https://doi.org/10.1016/S2215-0366(22)00312-1.

22. Parker RB, Larkin T, Cockburn J. Gender bias in medical images affects students' implicit but not explicit gender attitudes. *AERA Open*. 2018;4(3). https://doi.org/10.1177/2332858418798832.

23. Fingerhut AW, Martos AJ, Choi SK, Abdou CM. Healthcare stereotype threat and health outcomes among LGB individuals. *Ann Behav Med*. 2022;56(6):562-572. https://doi.org/10.1093/abm/kaab060.

24. OpenAI. Research Dall-e 2 pretraining mitigations. https://openai.com/research/dall-e-2-pre-training-mitigations. Accessed April 12, 2024.

25. Teresa-Morales C, Rodríguez-Pérez M, Araujo-Hernández M, Feria-Ramírez C. Current stereotypes associated with nursing and nursing professionals: an integrative review. *Int J Environ Res Public Health*. 2022;19(13):7640. https://doi.org/10.3390/ijerph19137640.

26. van Kolfschooten H. The AI cycle of health inequity and digital ageism: mitigating biases through the EU regulatory framework on medical devices. *J Law Biosci*. 2023;10(2):lsad031. https://doi.org/10.1093/jlb/lsad031.

27. Helberger N, Diakopoulos N. ChatGPT and the AI act. *Internet Policy Rev*. 2023;12(1). https://policyreview.info/essay/chatgpt-and-ai-act. Accessed May 29, 2023.

28. European Parliament. Generative AI and watermarking. *Briefing*. https://www.europarl.europa.eu/thinktank/en/document/EPRS_BRI(2023)757583. Accessed December 22, 2023.