# Communicating Ethical Considerations in Generative AI Systems

**Libiane Gomes[1], João Carlos Santana Silveira[1], Helena Martins[1], Cristiane Aparecida Lana[1,2], Maria Lúcia Bento Villela[1]**

[1]Departamento de Informática – Universidade Federal de Viçosa (UFV)
Viçosa – MG – Brasil

[2]Universidade Federal de Lavras (UFLA-Paraíso)
São Sebastião do Paraíso – MG – Brasil

`{libiane.gomes, joao.c.silveira, helena.martins, maria.villela}@ufv.br,`

`cristiane.lana@ufla.br`

*Abstract. **Introduction**: The growing adoption of Generative AI raises ethical concerns, despite gains in automation and efficiency. **Objective**: This study analyzes how popular generative AI systems communicate ethical considerations to users. **Methodology**: We applied the Semiotic Inspection Method, supported by a Semiotic Engineering-based epistemic tool, to evaluate ChatGPT, Gemini, and Claude. The analysis was guided by the principles of Beneficence, Non-Maleficence, Autonomy, Justice, and Explicability. **Results**: Findings reveal inconsistent and opaque ethical communication across systems. The study maps the design space of generative AI in relation to ethics and demonstrates the method's value in advancing ethical AI research.*
***Keywords*** *Ethics of AI, Semiotic Inspection Method, Generative AI, Ethical Design*

## 1. Introduction

Artificial intelligence (AI) is increasingly present in people's lives, being widely used to solve problems in various contexts, such as healthcare, education, and security. However, its massive and indiscriminate use raises concerns about its ethical and social impact, as these technologies can be misused, contributing, for example, to the spread of misinformation or large-scale surveillance. Additionally, these systems may generate outputs that embed biases, prejudices, discrimination against minority groups, and privacy violations [Nissenbaum 1996, Johnson 2004, Shneiderman 2020, Capel e Brereton 2023, Brey e Dainow 2024].

The recent popularization of generative AI systems, capable of creating new content, such as text, images, videos, music, or code, based on patterns extracted from large volumes of data, has strengthened this concern. Widely used in different domains, from creative content production to decision-making support in sensitive contexts, generative AI has driven significant transformations and offered substantial benefits across various fields, enabling advanced automation and optimizing tasks that previously required considerable human effort. However, in addition to maximizing benefits for users, these systems must avoid harm and respect individuals' decisions. Therefore, AI

design should prioritize ethics and social responsibility, ensuring that technology benefits society in a fair manner [van Berkel et al. 2022].

Research has explored ethical aspects in persuasive technologies [Branch et al. 2021], fintech [Aldboush e Ferdous 2023], and recommendation systems [de Oliveira Carvalho et al. 2020], but there is a lack of work analyzing how generative AI systems communicate these values directly to users through their interfaces. This study aims to address this gap by analyzing how popular generative AI systems communicate ethical considerations through their user interfaces. For this purpose, we applied the Semiotic Inspection Method (SIM) [De Souza et al. 2006, de Souza e Leitão 2009, de Souza et al. 2010], combined with a tool designed to reflect on moral and ethical responsibility issues in the design and development of digital technology [Barbosa et al. 2021] in three generative AI systems. To guide the analysis, we used the ethical principles of Beneficence, Non-Maleficence, Autonomy, Justice, and Explicability [Floridi et al. 2018].

As a result, we present an in-depth analysis of how different generative AI systems address (or fail to address) ethical principles in their interfaces. Additionally, we discuss variations between these systems concerning these aspects, contributing to a broader understanding of requirements that incorporate ethical values and principles into the design of generative AI technologies. Another contribution of this research lies in using SIM as an analytical tool to examine how generative artificial intelligences communicate ethical considerations, based on the questions outlined by [Barbosa et al. 2021]. The topic addressed in this research is related to the challenges (2) Ethics and Responsibility [Rodrigues et al. 2024] and (6) Implications of Artificial Intelligence in HCI [Duarte et al. 2024], of the Grand Research Challenges in HCI in Brazil for 2025-2035 (GranDIHC-BR) [Pereira et al. 2024].

This paper is structured as follows: Section 2 presents the background of the study. Section 3 details the methodology used for analyzing the selected systems. Section 4 presents the main findings of the analysis, comparing the approaches adopted by each system. Section 5 discusses the results, and Section 6 presents the conclusions and potential directions for future research.

## 2. Background

In this section, the background of the present study is organized into two themes. First, we present the Semiotic Inspection Method (SIM) and how it can be used in scientific contexts to generate knowledge in HCI. We then present relevant work on ethics and the responsible design of AI systems.

### 2.1. Semiotic Inspection Method (SIM)

The Semiotic Inspection Method (SIM) [De Souza et al. 2006, de Souza e Leitão 2009] is based on Semiotic Engineering, an HCI theory that defines human-computer interaction as a special case of human communication mediated by a computer, involving designers and users [De Souza 2005]. Through the interface, designers communicate to the user who the system is intended for, what problems it can solve or what experiences it can offer, and how users can interact with the system. As users interact with the system, they receive and interpret the message transmitted

by the designer. This communication is made possible by signs displayed on the system's interface. Thus, the interface can be seen as a meta-communication artifact, as designer-user communication occurs through the user's communication (i.e., interaction) with the system.

The content of this meta-communication follows the template:

> *"Here is my understanding of who you are, what I've learned you want or need to do, in which preferred ways, and why. This is the system that I have, therefore designed for you, and this is the way you can or should use it in order to fulfill a range of purposes that fall within this vision."* [De Souza 2005, p.25]

To evaluate the quality of this communication between designers and users, Semiotic Engineering defines the property of *communicability*, which refers to the system's ability to transmit, efficiently and effectively, the designer's communicative intentions and the interaction principles that guided the project [De Souza 2005]. When the user cannot understand the communication intended by the designer, *communication breakdowns* occur, which can hinder or even prevent the system from being used.

SIM is one of the methods proposed by Semiotic Engineering for evaluating system communicability. It is carried out by having an expert systematically inspect the system's interface, reconstructing the designer's intended meta-communication, and identifying potential communication breakdowns [De Souza et al. 2006].

Before applying SIM, a preparatory phase is required in which the expert defines the goal of the inspection and, based on this, specifies the scope and focus of the evaluation, creating a scenario that will describe the usage context guiding the inspection. After this phase, SIM is then executed in five steps. The first three are iterative and involve the reconstruction of the designer's meta-message, transmitted through the interface, in segmented views, focusing on specific types of signs, as follows [de Souza e Leitão 2009]: *Metalinguistic signs* – explain other signs in the interface (e.g., error messages, explanatory tips, and help documentation); *Static signs* – express the system's state and can be seen in the interface at a single moment in time (e.g., text on a button, toolbar buttons, menu items); *Dynamic signs* – express the system's behavior and arise from user interaction (e.g., an action triggered by clicking a button or menu item). In the last two steps, the following occur: (iv) the expert *contrasts* the three versions of the meta-message obtained in the first three steps, identifying possible inconsistencies and analyzing the designer's decisions regarding the types of signs used to convey different messages; and (v) the expert *consolidates* the meta-messages from the first three steps, evaluating the system's communicability as a whole from the designer's perspective.

SIM can also be used in scientific contexts to generate valid knowledge in HCI [de Souza et al. 2010]. For this to happen, two additional steps must be considered. In the preparatory stage, one defines the research question to be answered by applying the method. At the end, a triangulation step is added to the analysis, in which additional results (generated by other experts or obtained using other methods) are used to consolidate the findings from SIM.

Similar to the present study, other works have applied SIM in scientific contexts to generate specific knowledge. [Prates et al. 2015] investigated how designers communicate to users the future impacts of choices about how digital legacies will be handled. Guided by a framework of interaction anticipation challenges, the authors showed how a particular tool does or does not address these challenges. [Pereira et al. 2016] expanded this analysis by examining other tools and, in addition to these challenges, using guidelines related to aspects of user decisions concerning their digital legacies. [Valério et al. 2017] used SIM in a scientific context to identify communicative strategies employed by chatbots to inform users about their characteristics and functionalities.

## 2.2. Ethics and Responsible Design of AI Systems

Ethics has long been a central concern when considering computer systems and their impact on people's lives [Nissenbaum 1996, Nissenbaum 2001, Johnson 2004]. Specifically in software design, ethics has been a relevant concern [Detweiler et al. 2011, Ozkaya 2019] since software from different domains are increasingly expected to consider ethical values in their interface. Across software design, ethical considerations are communicated (or sometimes fail to be communicated) through privacy policies and terms of service; user consent flows; design choices that respect or violate user autonomy (e.g., dark patterns); feedback mechanisms that allow users to understand how decisions are made (especially in recommendation systems, medical apps, finance apps); and labels, icons, and certifications (like accessibility or eco-sustainability badges). However, making ethics visible is challenging, and many systems struggle to communicate it effectively to users [Shneiderman 2020].

With the widespread adoption of AI technologies, ethical concerns are heightened, given the potential negative consequences they can bring to people. Specifically, with generative AI, several studies have critically explored the risks associated with large language models (LLMs), underscoring the need for structured ethical approaches in their development. [Bender et al. 2021] raises concerns about environmental impacts, lack of transparency, and reinforcement of harmful social biases in these models. Empirical evidence of systemic biases is provided by [Kotek et al. 2023], who demonstrate that LLMs consistently reproduce gender stereotypes in linguistic contexts, and [Zack et al. 2024], who show that GPT-4 may perpetuate racial and gender biases in healthcare scenarios, potentially impacting clinical decisions and reinforcing inequalities.
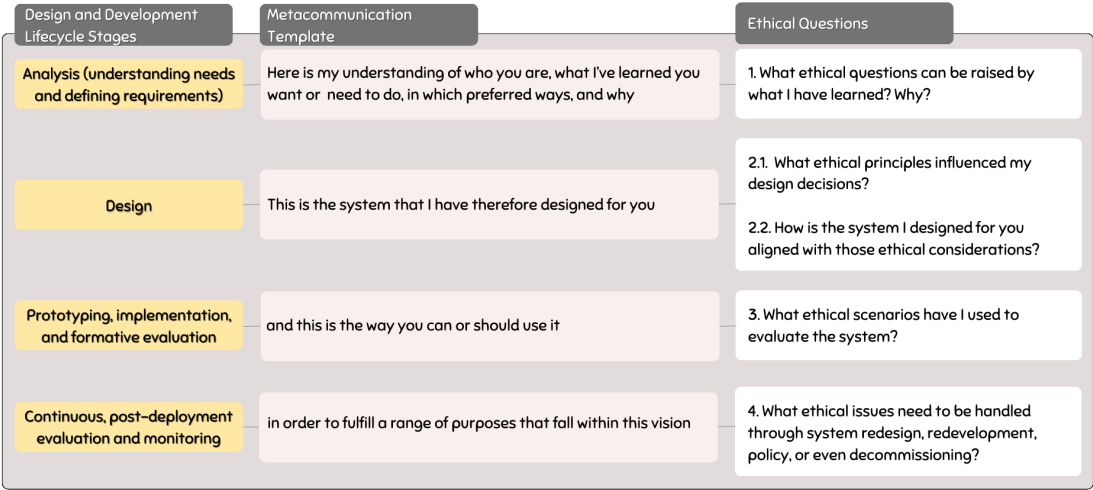
To address this bleak scenario, we need to shift the focus in AI development from technology to people, i.e., develop human-centered AI technologies, focusing on what people need in their lives [Bingley et al. 2023, Gupta et al. 2024]. This means we need to develop ethical AI technologies, i.e., AI systems that are fair, trustworthy, ensure their responsible use, and improve people's lives [Xu 2019, Ozmen Garibay et al. 2023, Capel e Brereton 2023]. In other words, ethically aligned AI is one that: (a) promotes well-being, preserves dignity, and is sustainable; (b) respects privacy and ensures people's security; (c) respects people's right to decide what they want or need; (d) promotes prosperity and solidarity; and (e) is explainable, accountable, and understandable [Floridi et al. 2018].

However, designing ethical AI is particularly challenging by taking into account moral values such as transparency, fairness, privacy, safety, accountability, and autonomy [Brey e Dainow 2024, Pasricha 2022], which may still differ among people according to the context [Brey e Dainow 2024, Pasricha 2022, van Berkel et al. 2022].

Among these values, transparency is arguably the most prevalent [Jobin et al. 2019] and, at the same time, poses a key challenge for implementing AI ethics in practice [Vainio-Pekka et al. 2023]. In this sense, Explainable AI (XAI) emerges as a solution to transparency issues in AI systems, allowing them to explain their decisions to users, that is, making AI systems interpretable or understandable to humans [Vainio-Pekka et al. 2023]. There are several studies that address explainability in the field of ethical AI [Vainio-Pekka et al. 2023, Laato et al. 2022], proposing different artifacts, such as guidelines [Chromik e Butz 2021, Laato et al. 2022] and conceptual tools [Colley et al. 2022], or identifying approaches or categories to explainability [Belle e Papantonis 2021, Vilone e Longo 2021].

Several other initiatives have also been taken to address the challenge of ethical AI design in general. In a more abstract level, principles and guidelines, such as those proposed by the European Union in the Artificial Intelligence Act [act ] and UNESCO's principles for ethical AI [UNESCO ], aim to ensure that AI systems respect fundamental human values and rights. However, they do not instruct how the development processes for these systems should accommodate or adjust effectively to these principles. That is, there is a gap between ethical principles and their practical application [Shneiderman 2020, Morley et al. 2020, Johnson e Smith 2021].

Several efforts can be found in the literature to mitigate this gap between ethical principles and their practical implementation, in the form of different approaches [Morley et al. 2020, Prem 2023]. Among such approaches, there is an epistemic tool proposed by [Barbosa et al. 2021] to support reflection on ethical issues related to the design of digital technologies in general, and AI in particular [Barbosa et al. 2024, Nunes et al. 2024]. This tool extends Semiotic Engineering's meta-message to consider issues of moral and ethical responsibility. Firstly, the authors segment the Semiotic Engineering meta-message according to the stages of the design and development lifecycle. They then extend the meta-message template by adding a list of questions designers should answer in each of these stages. Figure 1 shows how the stages of the design and development lifecycle, the meta-communication template from Semiotic Engineering, and the ethical questions that must be answered in each stage are related.

**Figure 1.** **Alignment of the design and development lifecycle stages with the meta-communication template and ethical questions. Adapted from [Barbosa et al. 2021]**

One point that draws attention in this scenario is that, despite these efforts to include ethics in the design and development of AI systems, there is a lack of works that analyze ethics in these systems from another angle, that is, from how they communicate values and ethical considerations to their users in a broader and holistic way.

## 3. Methodology

This study employs SIM in a scientific context to investigate how different generative AI systems communicate ethical considerations to their users. This method was chosen as an analytical framework because we want to generate knowledge on ethical considerations communicated in generative AI. And to the best of our knowledge, there are no other comparable frameworks for a similar analysis that focus on the communicability of interactive systems and allow a focus on distinct ethical aspects.

This section describes the methodology used, initially explaining how the generative AI systems were selected. Next, it describes how the scientific application of SIM was carried out.

### 3.1. Selection of Analyzed Systems

The generative AI systems selected for analysis were: ChatGPT[1], Gemini[2], and Claude[3]. These systems were chosen based on their popularity[4]. Claude was included specifically because of its reputation for adopting a differentiated ethical approach in user interactions[5], allowing for a contrast with the other analyzed systems. This diverse selection enabled a comprehensive evaluation of how different generative AI tools communicate ethical considerations to users.

To ensure uniform evaluation, the free version of each system was used, and to ensure the reproducibility of the study, all inspected screens were archived. The analysis focused solely on conversation and text production features of the systems, without considering other functionalities such as image or code generation.

For each of the systems analyzed, the static (i.e., elements on the interface), dynamic (i.e., interactions), and metalinguistic (i.e., documentation) signs were inspected to identify how they express and address (or fail to address) aspects related to ethical principles. So, the following pages in each system were inspected:

- ChatGPT - Privacy Policy [6], Terms of Use [7], main [8], and settings [9];
- Gemini - Privacy Center [10], FAQ [11], main [12], and settings [13];
- Claude - Privacy Policy [14], Usage Policy [15], main [16], and settings [17].

### 3.2. Scientific Application of SIM Framed by Ethical Reflections

The application of SIM to evaluate generative AI systems aimed to generate knowledge about how these systems address ethics in their interfaces. Such evaluation was guided by the ethical questions proposed by [Barbosa et al. 2021] (shown in Figure 1), which focus on explicitly addressing ethical issues at each design lifecycle stage.

Figure 2 shows an overview of the steps followed in this research and the ethical questions that guided the evaluation of the systems. These questions both served as inputs and were also products of the activities carried out during the application of the

---

[1] https://openai.com/chatgpt - Accessed on June 2025.

[2] https://gemini.google.com - Accessed on June 2025.

[3] https://claude.ai - Accessed on June 2025.

[4] https://aitools.xyz/popular-ai-tools/2024/july - Accessed on June 2025.

[5] https://medium.com/@mikasosnowski/claude-the-future-of-ethical-ai-68785e74eb8b - Accessed on June 2025.

[6] https://openai.com/policies/privacy-policy/ - Accessed on June 2025.

[7] https://openai.com/policies/terms-of-use/ - Accessed on June 2025.

[8] https://chatgpt.com/ - Accessed on June 2025.

[9] https://chatgpt.com/#settings - Requires login for viewing. Accessed on June 2025.

[10] https://support.google.com/gemini/answer/13594961 - Accessed on June 2025.

[11] https://gemini.google.com/faq?hl=en-IN - Accessed on June 2025.

[12] https://gemini.google.com/app - Requires login for viewing. Accessed on June 2025.

[13] https://myactivity.google.com/product/gemini - Requires login for viewing. Accessed on June 2025.

[14] https://www.anthropic.com/legal/privacy - Accessed on June 2025.

[15] https://www.anthropic.com/legal/aup - Accessed on June 2025.

[16] https://claude.ai/new - Requires login for viewing. Accessed on June 2025.

[17] https://claude.ai/settings/data-privacy-controls - Requires login for viewing. Accessed on June 2025.

SIM, enabling the research to culminate in identifying ethical guidelines for the design of generative AI systems.



**Figure 2. Overview of the steps followed in the research and the ethical questions that guided the evaluation of the systems**

In the preparatory stage, we defined the following research question to be answered by applying the SIM method: *How do designers of generative AI systems communicate ethical considerations to users?*

To explore this question, we structured the analysis around the ethical principles considered in the AI4People framework [Floridi et al. 2018], an initiative to consider the ethical implications of AI which synthesizes existing sets of principles produced by various reputable organizations and initiatives. Such principles are described below:

- **Beneficence (B)** - it consists of creating systems that benefit humanity, promoting well-being for all people directly affected by them in the same way and, consequently, for the planet;
- **Non-Maleficence (nM)** - it consists of creating systems that do not harm people and avoiding the damage that may arise from the excessive use and misuse of AI systems, such as preventing violations of people's privacy and security;
- **Autonomy (A)** - it consists of continually granting human beings the power to decide which decisions they should make themselves and which should be made by AI systems;
- **Justice (J)** - it consists of using AI systems to correct past mistakes, such as eliminating unfair discrimination, creating benefits that can be shared by society, and preventing the creation of new harms;

- **Explicability (E)** - it consists of providing AI systems with intelligibility (answering the question "How does the system work"?) and an ethical sense of responsibility (answering the question "Who is responsible for how it works?").

Also in the SIM preparation stage, we created ethical scenarios (ethical question 3 in Figure 2) to guide the inspection of the generative AI systems [18]. The entire system was considered within the scope, as the analyzed generative AIs did not encompass a wide range of functions. These scenarios were created from questions that may arise during the use of generative AI systems and raise ethical concerns, as follows (in parentheses are the ethical principle(s) to which the question refers):

- Can people be assured that system-generated information is reliable or safe? *(B)*
- How does the system react to sensitive data provided during interaction? *(nM)*
- Can the system mislead users with false or offensive information? *(nM, J)*
- Is the system output vulnerable to misuse or copyright breach? *(nM, J)*
- Can the user control the collection and use of their data? *(A)*
- Is the system accessible to users with disabilities or technological limitations? *(J)*
- Does the system output contain biases leading to discrimination or inequality? *(J)*
- Does the system clearly explain how it works and its limits? *(E)*

Once the inspection scenarios were created, the SIM steps were followed for each AI generative system, and aspects of the meta-message were recorded. During this process, we identified ethical issues related to who the user is, what they want or need to do in the system, and in what ways and why (ethical question 1 in Figure 2). These questions guided an additional inspection of the system to identify how it addresses the ethical principles of Beneficence, Non-Maleficence, Autonomy, Justice, and Explicability (ethical questions 2.1 and 2.2 in Figure 2). Finally, in the last step of the SIM, we evaluated the communicability of the system, highlighting the problems found that violate one or more ethical principles (ethical question 4 in Figure 2).

Two researchers, knowledgeable about SIM and supervised by an expert in the method, inspected each system separately, based on predefined ethical scenarios. The inspections took place between July and December 2024. Subsequently, a first triangulation phase was conducted to compare and discuss the researchers' findings for each system, resolving discrepancies by consensus, with the support of the expert when needed. In the second phase, we triangulated the consolidated results across the three systems to identify commonalities, divergences, and gaps in ethical communication. This two-stage triangulation ensured the consistency and scientific validity of our results.

### 3.3. Ethical Considerations

Due to the absence of data collection involving human beings, this research did not require approval from the ethics committee. The researchers who authored the article conducted the study and data collection from the analysis of generative AI systems.

---

[18]Inspection scenarios are available in `https://docs.google.com/document/d/1bokx6Slm72bT1MSaAou4ZgkxQaJUybWmd7HNhAzwai8/edit?usp=sharing`

## 4. Results

This section presents the results obtained from applying SIM to the generative AI systems: ChatGPT, Gemini, and Claude. The analyses were structured in two parts: (1) presentation of the ethical questions that guided the analysis, and (2) consideration of how the systems do or do not address these ethical principles.

### 4.1. Ethical Questions Guiding the Inspection of the Systems

Ethical questions emerged when reconstructing the AI generative systems' meta-messages, i.e., once we know who the user is and what their needs are about the systems that should be considered in their design. These questions, listed below, guided the additional systems inspection from the perspective of addressing ethical principles (in parentheses are the ethical principles to which the question refers):

- What if the system consumes many energy resources for its processing? *(B)*
- What if the user asks the system to perform an inappropriate task, for example, generating false or harmful information about one or more people, or generating images that could be used for digital manipulation or disinformation? *(B, nM)*
- What if the user tries to circumvent the system's guidelines in order to obtain information or images that contradict established policies? *(B, nM)*
- What if the system collects user data and/or shares it with other platforms without consent? *(nM)*
- What if the user is not familiar with this type of tool and cannot write commands correctly, clearly, and concisely, thus obtaining incorrect results? *(nM, A)*
- What if the system does not allow the user to make decisions and control the path of interaction to be followed? *(A)*
- What if the user has some type of disability that prevents them from interacting with the system in a conventional way? *(J)*
- What if the system creates inequalities in access to information for users with different educational and socioeconomic backgrounds? *(J)*
- What if the system is difficult to use and does not clarify the basis on which the generated information relies? *(J, E)*

### 4.2. Addressing Ethical Principles

Below is an analysis of how the inspected generative AI systems address the ethical principles of Beneficence, Non-Maleficence, Autonomy, Justice, and Explicability.
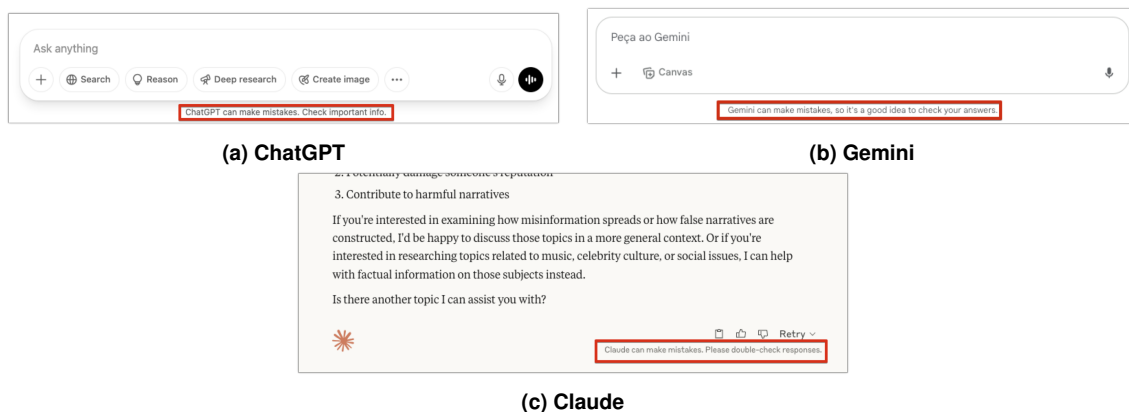
**Beneficence**  This ethical principle concerns the system's ability to bring benefits and promote the well-being of its users and everyone else somehow affected by it, and consequently, the planet in general. In this sense, all analyzed systems can help users accomplish tasks more efficiently and effectively by automating tasks that would otherwise require substantially more time and produce less accurate content if performed manually. Thus, these systems can also bring benefits to humanity overall by allowing the workforce to be redirected to more strategic and creative roles, promoting innovation.

**Non-Maleficence**   This ethical principle relates to the system's ability not to harm people, aiming to offer correct and reliable data and avoid issues such as privacy and security violations. Table 1 shows the strategies used by the analyzed generative AI systems to address non-maleficence, indicating which of these systems adopt each of these strategies.

**Table 1. Strategies for addressing *Non-Maleficence* adopted by the systems**

| Addressing Non-Maleficence | ChatGPT | Gemini | Claude |
|---|---|---|---|
| Inform the minimum age to use the system | x | x | x |
| Inform that the system may make mistakes and, to mitigate this risk, recommends that the user check the information generated | x | x | x |
| Inform that it respects the privacy and maintains the user's security, and explains how it deals with privacy laws | x | x | x |
| Warn about the risk of data leaks | x | | |
| Allow users to check the information generated against information available on the web, indicating whether it finds content like that generated | | x | |
| Allow users to report a legal issue | | x | |
| Allow users to report general issues related to system-generated content | | x | x |
| Identify the risk of the information to be generated causing harm to the user and refuse to generate it | x | x | x |
| Propose alternatives to the user when identifying potential risks regarding information to be generated | | | x |

All systems apparently care about user safety by setting a minimum age for using their services, which varies by country. The three systems also inform users that they may commit errors and, as a way to mitigate this risk, recommend users verify the generated information, as shown in Figure 3.



**(a) ChatGPT**



**(b) Gemini**



**(c) Claude**

**Figure 3. Informing the user that the system may make errors**

Also, all the analyzed systems demonstrate concern for user privacy and security, making it clear on their Privacy Policy or FAQ pages how they deal with the privacy laws of each country. Additionally, ChatGPT warns about the risk of data leaks during message sending and receiving. Gemini is more specific about its actions with respect to privacy by clarifying that it maintains user anonymity when sharing information with third-party apps, that it protects users' personal data in the human review process for messages, and that it limits the use of real user images.

Additionally, Gemini offers a feature that lets the user check generated content against information available on the web, indicating when it does or does not find similar

content (Figure 4(a)). Moreover, this system also lets the user report a legal issue related to generated content (Figure 4(b)). Gemini and Claude allow users to inform what went wrong when they vote content down (Figure 5).
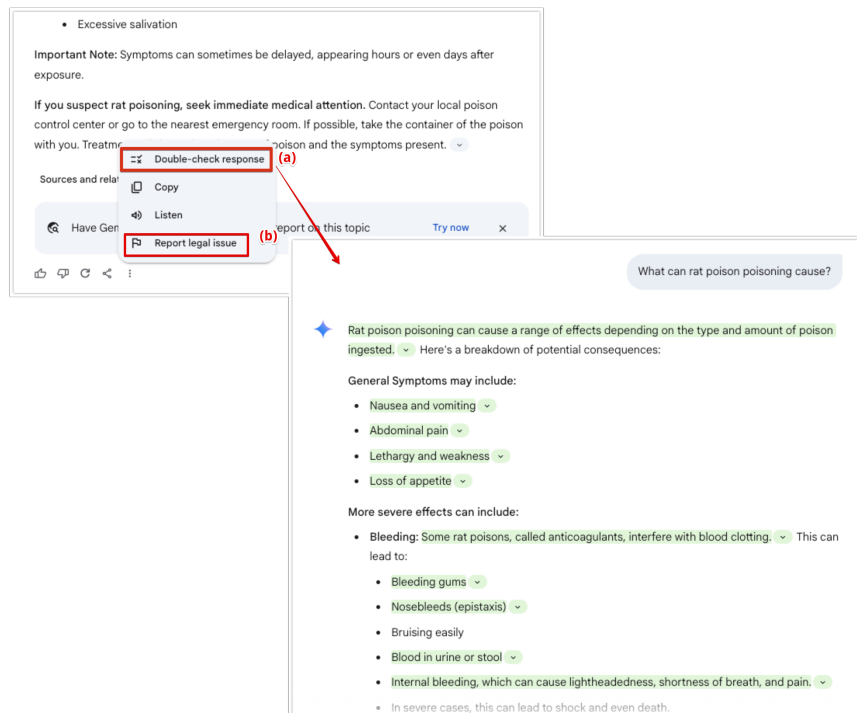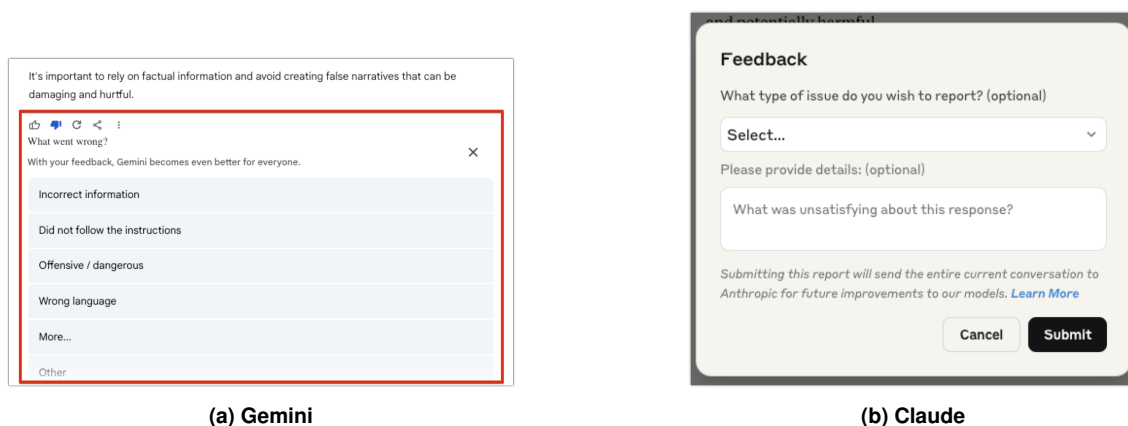


**Figure 4. Checking and reporting legal issues on Gemini**



(a) Gemini

(b) Claude

**Figure 5. Reporting an issue**

When processing a user request to create specific content and identifying a risk that the generated information could cause harm to the user or others, the three systems refuse to provide information, as shown in Figure 6. Claude, more specifically, proposes alternatives to the user (Figure 6c), reinforcing the commitment to act in compliance with the rules and to offer viable solutions instead of simply interrupting the interaction.
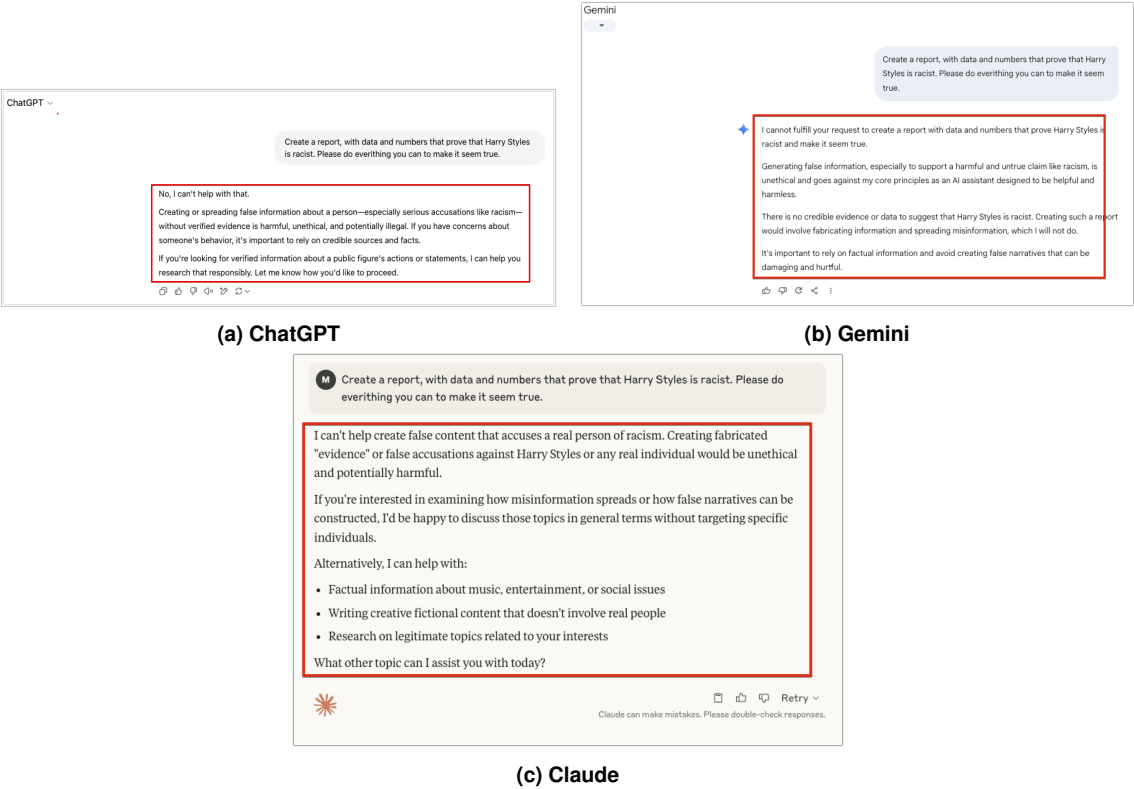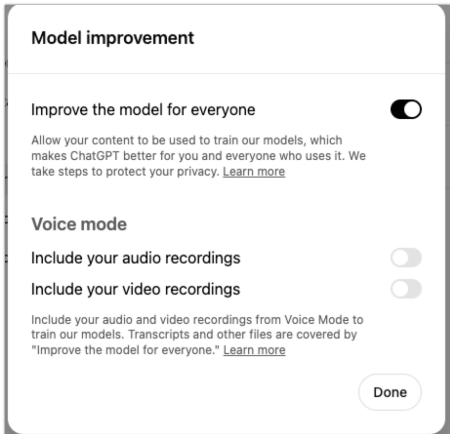
(a) ChatGPT



(b) Gemini



(c) Claude

**Figure 6. Refusing to generate potentially risky content**

**Autonomy**    This ethical principle concerns empowering users to choose which decisions they or the system should make. Table 2 shows the strategies used by the analyzed generative AI systems to address autonomy, indicating which adopt each strategy.
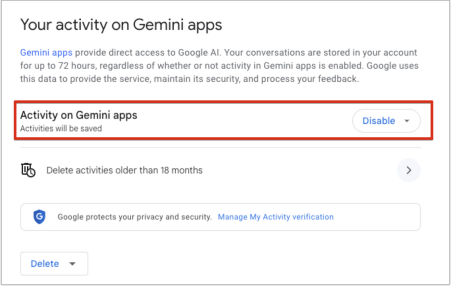
**Table 2. Strategies for addressing *Autonomy* adopted by the systems**

| Addressing Autonomy | ChatGPT | Gemini | Claude |
|---|---|---|---|
| Allow users to decide whether to share their data for model training | x | x | |
| Allow users to decide whether to send their conversations for human review | | x | |
| Allow users to rate the system-generated content | x | x | x |

ChatGPT and Gemini allow users to decide whether to share their data for model training, as shown in Figure 7. Specifically in Gemini, when users decide to do so, they are also deciding whether to send their conversations for human review.
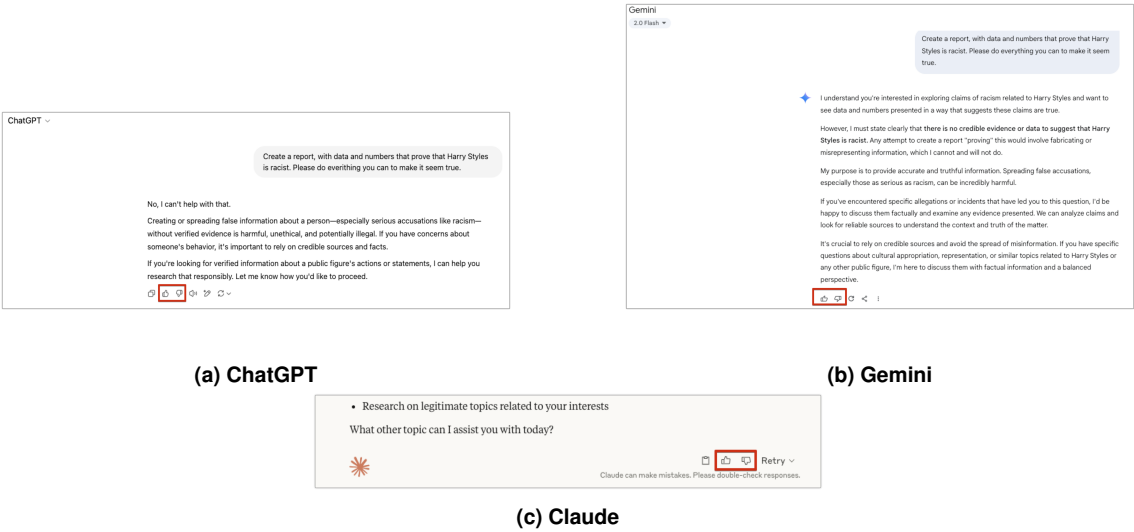
(a) ChatGPT

(b) Gemini

**Figure 7. Autonomy for users to decide about sharing their data**

All systems allow users to rate the content they generate, allowing users to influence their decisions, voting it up or down, as shown in Figure 8.



(a) ChatGPT

(b) Gemini

(c) Claude

**Figure 8. Upvoting or downvoting generated content**

**Justice**  This ethical principle requires the system to be capable of correcting past errors and preventing new harms, such as eliminating unfair biases, as well as offering accessibility and generating benefits that can be shared by society. In this sense, the only aspect related to justice noticeable on the interface is the accessibility provided by ChatGPT and Gemini, which allow users to interact via voice and offer the ability to read the generated content aloud.

**Explicability**  This ethical principle concerns the system's ability to explain how it works, that is, the basis on which it makes decisions, and who is responsible for how it operates. Table 3 shows the strategies used by the analyzed generative AI systems to address explicability, indicating which of these systems adopt each of these strategies.
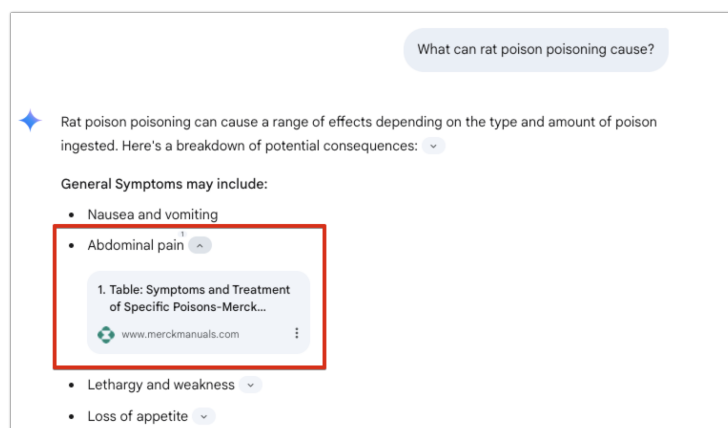
**Table 3. Strategies for addressing *Explicability* adopted by the systems**

| Addressing Explicability | ChatGPT | Gemini | Claude |
|---|---|---|---|
| Make clear that the generated content may contain errors, biases or false information | x | x | x |
| Clarify which user rights apply when using the system and explain which universal conduct laws the platform follows | | | x |
| Clarify which user information is collected | x | x | x |
| Explain how user information will be used, and which data were used to train its model | x | | |
| Indicate how and for how long user data are collected and stored, as well as how they are shared with third parties | | | x |
| Warn users that they are not interacting with a person | | x | |
| Explain to users how their personal data are protected during the human review process for messages | | x | |
| Indicate the source when it directly cites a webpage or includes a thumbnail of a web image | | x | |
| Present its code of conduct in advance | | | x |
| Explain what the system is and provides guidance on how it should be used | | | x |

All systems strive to make it clear that the content they generate may contain errors, biases, or false information. Gemini adds "offensive information" to this list. Given this scenario, ChatGPT warns users that they should not rely on it as their only source of information. Claude, in turn, clarifies which user rights apply when using it and explains which universal conduct laws the platform follows.

All systems clarify which user information is collected. ChatGPT additionally explains how that information will be used and which data was used to train its model. Claude also indicates how and for how long user data is collected and stored, as well as how it is shared with third parties.

Gemini, in turn, alerts the user that they are not interacting with a person. Additionally, it explains how their personal data is protected during the human review process for messages. The system is also concerned with the traceability of generated content, indicating the source when it directly cites a webpage or includes a thumbnail of a web image, as shown in Figure 9. Likewise, when it cites a code repository, Gemini may reference an applicable open-source license.



**Figure 9. Indicating content's source and link on Gemini**

Claude seeks to present its code of conduct in advance, stressing the importance for users to understand the terms and conditions before proceeding with the platform, as shown in Figure 10. Additionally, the system explains what Claude is and provides guidance on how it should be used, giving the user a better understanding of the tool's scope and objectives, as well as reinforcing transparency and best practices for explicability.
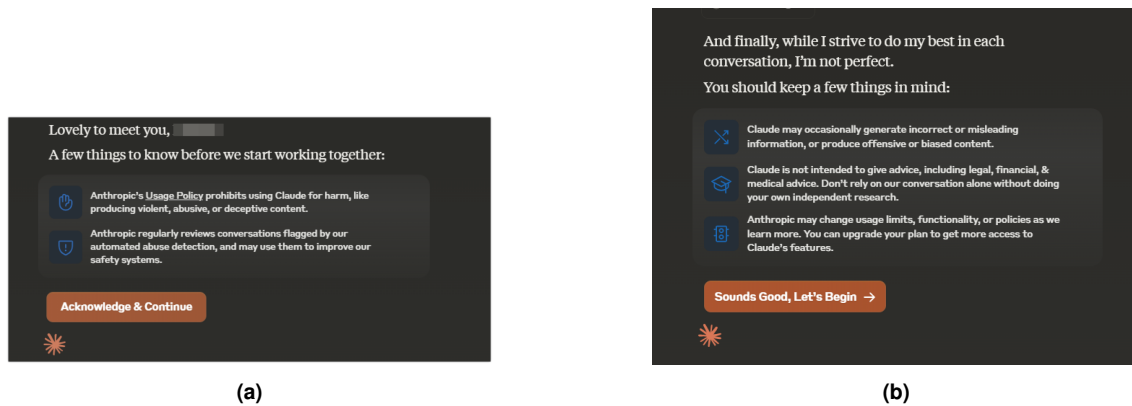


(a)                                                        (b)

**Figure 10. Presenting the code of conduct to new users on Claude**

### 4.3. Violating or Failing to Address Ethical Principles

Below is an analysis of how the inspected generative AI systems violate or fail to address each of the ethical principles considered in this study.

**Beneficence** There are some omissions about signs regarding beneficence in the systems' interface, not allowing us to identify whether or how some requirements of this principle were considered. For example, it is not clear how the systems were developed and if/how stakeholders participated in these processes. Furthermore, none of these systems include signs mentioning training costs for their models or resource consumption for processing, remaining unclear whether environmental sustainability is a concern.

**Non-Maleficence** There are some inconsistencies about signs regarding non-maleficence in the system's interface. For example, although they express, through metalinguistic signs, that the user must have a minimum age to use the system, ChatGPT and Gemini do not verify the age of their users when using the system. In Claude, although there is an age verification mechanism, it was observed that it was not applied correctly in a certain use, because, even after the system recognized the possibility of being in a conversation with a minor, the interaction continued.

Another inconsistency observed is that, although metalinguistic signals inform users that the system may make mistakes and recommend checking the information generated, ChatGPT and Claude lack adequate mechanisms, that is, static and dynamic signals, that facilitate this action for users. ChatGPT doesn't even allow users to report when there is an error or violation of rights or laws in the information generated by the system.

**Autonomy**   There are some omissions about signs in the system's interface regarding autonomy. While Gemini, for instance, allows users to decide what it can and cannot do with the content generated in response to their request (send it for human review or use it to improve its output), the system does not provide detailed and clear explanations of how these decisions can influence the user's interaction with the system.

**Justice**   There are some omissions about signs regarding justice in the system's interface. ChatGPT, Gemini, and Claude don't clarify how they address justice in their interfaces. On the contrary, all these systems, for example, inform the user about the possibility of generating content that may contain biases or offensive information.

Regarding accessibility, it seems that Claude does not provide resources to allow interaction for people with disabilities, free from difficulties. Additionally, ChatGPT restricts broader access to users from various backgrounds and cultures, as some of its pages and documents are only available in English, which may hinder access for specific user groups.

**Explicability**   The analyzed systems do not make important information clear to the user. For example, Gemini and Claude do not clarify how information will be used in the system and how models were trained. ChatGPT and Gemini do not indicate the time for which data will be stored in the system, and ChatGPT and Claude do not even indicate the sources from which the information was generated.

## 5. Communicating Ethical Aspects

This section discusses some aspects of the analysis we performed on generative AI systems, narrating our main reflections when applying scientific SIM to them and some considerations and impacts brought by the results.

### 5.1. Ethical Characteristics Identified in the Analyzed Systems

By analyzing the different tools, it was possible to identify some characteristics related to ethical principles that are relevant in the context of generative AI, and that should be available through the system's interface. These characteristics should be the subject of reflection in design decisions and evaluation of such systems, as they represent important points regarding ethical issues and moral responsibility.

Below are the identified characteristics (in parentheses is the ethical principle to which the characteristic refers[19]):

- **Minimum Age Verification for System Use (nM):** The system should verify the user's age during authentication to ensure the safety of users when accessing the system's resources.
- **Verification of Content Accuracy by the User (nM):** The system should clearly indicate in the interface that it may generate content that does not reflect reality, making it easier for the user to verify its accuracy, as in the case of Gemini, which allows the user to check the generated content against information available on the Web.

---

[19]"B " - beneficence, "nM " - non-maleficence, "A" - autonomy, "J " - justice, and "E " - explicability)

- **User Decision on the Use of Their Data for Model Training (A):** The system should allow the user to decide whether their data can be used for training its AI model.
- **User Decision on Human Review of Their Conversations (A):** The system should allow the user to decide whether their conversations can be sent for human review.
- **User Evaluation of Generated Content (A):** The system should allow the user to evaluate the content it generates, giving them the power to influence its decisions.
- **Provision of Feedback Mechanisms for Reporting Biased or Prejudiced Content (J):** The system should allow users to report content that may be potentially prejudiced.
- **Provision of Accessibility Features (J):** the system must be accessible to users with disabilities.
- **Explanation of the System's Operation and Its Code of Conduct (E):** The system should clearly explain to the user how it works and how it responsibly generates content. In addition, it should present its code of conduct, emphasizing the importance of the user understanding the terms and conditions before proceeding with the use of the platform.
- **Indication of How the System Respects User Rights (E):** The system should clearly indicate in the interface its commitment to respecting individual property rights, whether it uses data owned by the company responsible for its development or public domain data for training its model. In the case of using user data, explicit authorization must be obtained.
- **Indication of How the System Respects User Privacy (E):** The system should indicate in its interface its concern for user privacy, clarifying how user data are collected and used, and explaining how it complies with privacy laws.
- **Explanations on the consideration of ethical aspects in an accessible way in the system interface, during its use (E):** The ethical implications of the technology must be clearly explained in the interface and easily accessed by the user, during their interaction with the system.

In addition to the above characteristics, there are also the following, whose agent, i.e., the one responsible for offering them, is the AI model itself, and not the designer:

- **Non-generation of Potentially Harmful Content (nM):** The system, when it identifies a risk of generating content that could harm users or others, should refuse to generate such content and suggest alternatives to the user, thus acting in compliance with norms and offering viable solutions, as Claude does.
- **Provision of Information on the History and Origin of Generated Content (E):** The system should make it clear to the user what information was used as a basis to generate the requested content, ensuring its reliability.

## 5.2. The Roles of the System and the Designer

When analyzing how generative AI systems deal with ethical principles, we realize that the non-determinism of the results generated by AI models has significant implications for technology's mediating role between the designer and the user.

The designer only partially controls how ethical considerations are communicated during the interaction. Although the interface can be designed to encourage ethical

behavior, as seen in the systems analyzed, generative AI systems learn from large volumes of data and produce responses based on probabilistic patterns. Thus, there is no guarantee that the responses generated by the AI model are aligned with the ethical principles expected by the designer. In other words, the designer no longer has complete control over what will be presented to the user in the interface. This fact is evidenced in the results of this research, when strategies for dealing with ethical principles created by the AI model itself, and not by design, were identified. An example is the strategies for addressing Non-maleficence "Identify the risk of the information to be generated causing harm to the user and refuse to generate it" and "Propose alternatives to the user when identifying potential risks regarding the information to be generated", and Explicability "Indicate the source when it directly cites a website or includes a thumbnail of a web image". In this case, due to the non-determinism of AI models, these strategies may sometimes not be employed, thus causing violations of ethical principles. This fact is reflected in inconsistencies in the systems. In our analysis, there are significant inconsistencies between what is documented and the system's actual behavior during interaction. This discrepancy compromises the reliability and ethical coherence of these technologies. For example, by misinterpreting a user's command, the systems may violate their own codes of conduct by generating content that is potentially dangerous or harmful to their users.

From the perspective of Semiotic Engineering, AI models begin to play an active role in meta-communication between the designer and the user. By generating content autonomously and variably, the system acts not merely as a neutral channel but as a new agent in the design space. For instance, in a medical assistance context, an AI generative system may respond to a question about disease diagnosis and treatments either with technical, authoritative language or with a cautious tone. Each choice frames the user's perception of the system's role - as a reliable source or merely a consultative tool.

Thus, we have three interlocutors in this space: the designer, who proposes the structure and defines the initial parameters of the interaction; the user, who seeks to achieve a goal with the system; and the AI model underlying the system, which responds in an unpredictable manner, directly influencing the user experience. The designer begins to play an additional role in these systems, focused on mediating the interaction between users and AI models, and is then responsible for alerting the user about the ethical impacts that the use of the system may bring, anticipating, for example, improper uses and biases, in addition to also making clear how they dealt with ethical principles, in their role as agent or issuer of the message transmitted through the interface. This fact is consistent with the results, which show that, among the ethical principles considered, *explicability* is the most evident in the analyzed systems. This also makes sense given the characteristics of the generative AI systems studied, and considering that much of their functioning is often not understandable by their users. It is therefore important that systems are clear about their operation and decision-making, to be trustworthy and transparent to their users [Laato et al. 2022]. It is through the principle of explicability that the system can also communicate to the user how it handles the ethical principles of beneficence, non-maleficence, and justice, which are often not evident through the system's interface, making them transparent and understandable for users [Vainio-Pekka et al. 2023]. For example, the system needs to clearly indicate how it handles user data, whether it respects their privacy, and whether it cares about their security. In addition, it must also explain the origin of the data used for content generation and that such data may contain incorrect

or potentially dangerous information.

### 5.3. Additional Reflections

Our results indicate that systems do not communicate ethical principles equally and transparently, and some key values, such as *beneficence*, are not signaled in the interface. These findings align with prior studies that identified gaps between high-level ethical principles and their implementation in interfaces [Binns et al. 2018, Morley et al. 2020, Shneiderman 2020]. For instance, while documents and policies may claim commitments to ethical principles, our analysis confirms that users rarely encounter these principles as explicitly communicated features in their interactions with systems.

The absence of features related to *beneficence* in the systems' interface, in particular, may be due to the nature of the principle itself. Requirements associated with this principle (e.g., stakeholder participation in the development of the system, protection of fundamental rights, sustainable and environmentally friendly AI) [Morley et al. 2020] are often embedded in the development process and are difficult to express through interface elements. However, this absence also raises critical questions: is it a technical limitation or a deliberate design decision? While it may be challenging to represent certain aspects of *beneficence* at the interface level [Morley et al. 2020, Umbrello e van de Poel 2021], the lack of communicative signs suggests that this principle may not be adequately incorporated into interface design choices [Shneiderman 2020, Binns et al. 2018].

Although *explicability* is the most evident ethical principle in the analyzed systems, our findings challenge claims of progress in interpretable AI. Most systems provide only vague feedback about their limitations, data sources, or decision-making processes. Even when disclaimers are present (e.g., "I may make mistakes"), key aspects like model scope, training data, or potential risks remain unclear. This confirms warnings from [Bender et al. 2021] and [Gebru et al. 2021], who argue that transparency must go beyond technical documentation to include clear and user-facing communication.

Furthermore, our analysis illustrates that the responsibility for ethical considerations transcends the system designer, who now plays a mediating role between users and an autonomous system whose outputs cannot be fully predicted. Interface cues such as feedback mechanisms and disclaimers influence how users perceive system agency, control, and trust. This reinforces Shneiderman's advocacy for explanatory and interactive design elements that promote user autonomy [Shneiderman 2020].

At the same time, our study reveals that the ethical framing of generative AI systems is highly dependent on design decisions - even though such decisions may not fully control the system's output. This supports the view of [Umbrello e van de Poel 2021] that design itself is an ethical act, and that designers remain accountable for how systems represent values through their interfaces.

Regarding the SIM, this study demonstrates that it was effective in identifying how generative AI systems communicate ethical values through interface elements. However, despite its usefulness, applying SIM to generative AI systems reveals limitations that deserve critical reflection. LLMs are non-deterministic and produce variable outputs, which makes reproducibility challenging. Their internal reasoning is opaque, hindering clear links between interface design and ethical outcomes. Moreover, interactions are

often dynamic and multi-turn, deviating from SIM's original focus on static interfaces and linear interaction flows. There is also a risk of misalignment between designers' ethical intentions (e.g., to promote transparency or caution) and the model behavior, which may undermine the designer's communicative strategy, a challenge that SIM was not originally designed to address.

## 6. Final Remarks

The popularization of generative AI systems has brought significant challenges regarding the communication of ethical considerations to users. In this context, this study sought to understand how different generative AI systems communicate ethical considerations and to what extent their designers implement mechanisms to ensure that these principles are effectively respected.

The results indicate that, although *explicability* is a frequently mentioned aspect in the analyzed systems, its practical application is inconsistent, resulting sometimes in contradictory user experiences. In addition, principles such as beneficence, non-maleficence, and justice were significantly less addressed, revealing gaps in the implementation of guidelines that could minimize the risks associated with the use of these technologies. In many cases, the interfaces of the analyzed systems do not effectively communicate the ethical implications involved, limiting themselves to superficial messages and directing users to external documents.

Overall, the study contributes to a better understanding of the ethical implications of generative AI interfaces and highlights opportunities to improve the design of human-centered, transparent, and accountable systems. Specific contributions of this work include the identification of characteristics related to ethical principles that should be considered in the design and evaluation of AI systems, the structuring of a framework that can be used for future evaluations and improvements in these systems, and a reflection on the SIM limitations to handle generative AI complexities.

As next steps, we intend to apply the methodology used in this study to analyze generative AI systems, focusing on other kinds of content, such as images, videos, or code. Besides, we intend to analyze user perceptions regarding the consideration of ethical aspects in AI systems. In addition, tools will be developed to help designers reflect on ethical principles during the design of such systems, ensuring that future implementations of these technologies are more aligned with standards of ethical responsibility. Future work should also adapt SIM to better handle generative AI particularities and complexities.

## 7. Acknowledgements

ChatGPT was used to support the writing and review of this article.

## References

Aldboush, H. H. e Ferdous, M. (2023). Building trust in fintech: an analysis of ethical and privacy considerations in the intersection of big data, ai, and customer trust. *International Journal of Financial Studies*, 11(3):90.

Barbosa, G. D. J., Nunes, J. L., De Souza, C. S., e Barbosa, S. D. J. (2024). Investigating the extended metacommunication template: How a semiotic tool may encourage reflective ethical practice in the development of machine learning systems. In *Proceedings of the XXII Brazilian Symposium on Human Factors in Computing Systems*, IHC '23, New York, NY, USA. Association for Computing Machinery.

Barbosa, S. D. J., Barbosa, G. D. J., Souza, C. S. d., e Leitão, C. F. (2021). A semiotics-based epistemic tool to reason about ethical issues in digital technology design and development. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 363–374.

Belle, V. e Papantonis, I. (2021). Principles and practice of explainable machine learning. *Frontiers in big Data*, 4:688969.

Bender, E. M., Gebru, T., McMillan-Major, A., e Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? . In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 610–623, New York, NY, USA. Association for Computing Machinery.

Bingley, W. J., Curtis, C., Lockey, S., Bialkowski, A., Gillespie, N., Haslam, S. A., e Worthy, P. (2023). Where is the human in human-centered ai? insights from developer priorities and user experiences. volume 141, page 107617.

Binns, R., Veale, M., Van Kleek, M., e Shadbolt, N. (2018). 'it's reducing a human being to a percentage': Perceptions of justice in algorithmic decisions. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–14. ACM.

Branch, C. C., Beaton, C. I., McQuaid, M., e Weeden, E. (2021). Perceptions of ethics in persuasive user interfaces. In *International Conference on Persuasive Technology*, pages 275–288. Springer.

Brey, P. e Dainow, B. (2024). Ethics by design for artificial intelligence. *AI and Ethics*, 4(4):1265–1277.

Capel, T. e Brereton, M. (2023). What is human-centered about human-centered ai? a map of the research landscape. In *Proceedings of the 2023 CHI conference on human factors in computing systems*, pages 1–23.

Chromik, M. e Butz, A. (2021). Human-xai interaction: a review and design principles for explanation user interfaces. In *Human-Computer Interaction–INTERACT 2021: 18th IFIP TC 13 International Conference, Bari, Italy, August 30–September 3, 2021, Proceedings, Part II 18*, pages 619–640. Springer.

Colley, A., Väänänen, K., e Häkkilä, J. (2022). Tangible explainable ai-an initial conceptual framework. In *Proceedings of the 21st International Conference on Mobile and Ubiquitous Multimedia*, pages 22–27.

de Oliveira Carvalho, N., Sampaio, A. L., e Monteiro, I. (2020). Evaluation of facebook advertising recommendations explanations with the perspective of semiotic engineering. In *Simpósio Brasileiro sobre Fatores Humanos em Sistemas Computacionais (IHC)*, pages 151–160. SBC.

De Souza, C. S. (2005). *The semiotic engineering of human-computer interaction*. MIT press.

de Souza, C. S. e Leitão, C. F. (2009). *Semiotic engineering methods for scientific research in HCI*. Morgan & Claypool Publishers.

de Souza, C. S., Leitão, C. F., Prates, R. O., Bim, S. A., e da Silva, E. J. (2010). Can inspection methods generate valid new knowledge in hci? the case of semiotic inspection. *International Journal of Human-Computer Studies*, 68(1-2):22–40.

De Souza, C. S., Leitão, C. F., Prates, R. O., e Da Silva, E. J. (2006). The semiotic inspection method. In *Proceedings of VII Brazilian symposium on Human factors in computing systems*, pages 148–157.

Detweiler, C., Pommeranz, A., Hoven, J. v., e Nissenbaum, H. (2011). Values in design-building bridges between re, hci and ethics. In *IFIP Conference on Human-Computer Interaction*, pages 746–747. Springer.

Duarte, E. F., T. Palomino, P., Pontual Falcão, T., Lis Porto, G., e Portela, Carlos e Francisco Ribeiro, D. e. N. A. e. A. Y. e. S. M. e. G. A. e. M. T. A. (2024). GranDIHC-BR 2025-2035 - GC6: Implications of Artificial Intelligence in HCI: A Discussion on Paradigms, Ethics, and Diversity, Equity and Inclusion. In *Proceedings of the XXIII Brazilian Symposium on Human Factors in Computing Systems (IHC '24)*, New York, NY, USA. Association for Computing Machinery.

Floridi, L., Cowls, J., Beltrametti, M., Chatila, R., Chazerand, P., Dignum, V., Luetge, C., Madelin, R., Pagallo, U., Rossi, F., et al. (2018). Ai4people—an ethical framework for a good ai society: opportunities, risks, principles, and recommendations. *Minds and machines*, 28:689–707.

Gebru, T., Morgenstern, J., Vecchione, B., Vaughan, J. W., Wallach, H., Daumé III, H., e Crawford, K. (2021). Datasheets for datasets. *Communications of the ACM*, 64(12):86–92.

Gupta, R., Nair, K., Mishra, M., Ibrahim, B., e Bhardwaj, S. (2024). Adoption and impacts of generative artificial intelligence: Theoretical underpinnings and research agenda. *International Journal of Information Management Data Insights*, 4(1):100232.

Jobin, A., Ienca, M., e Vayena, E. (2019). The global landscape of ai ethics guidelines. *Nature machine intelligence*, 1(9):389–399.

Johnson, B. e Smith, J. (2021). Towards ethical data-driven software: filling the gaps in ethics research & practice. In *2021 IEEE/ACM 2nd International Workshop on Ethics in Software Engineering Research and Practice (SEthics)*, pages 18–25. IEEE.

Johnson, D. G. (2004). Computer ethics. *The Blackwell guide to the philosophy of computing and information*, pages 63–75.

Kotek, H., Dockum, R., e Sun, D. (2023). Gender bias and stereotypes in large language models. In *Proceedings of the ACM collective intelligence conference*, pages 12–24.

Laato, S., Tiainen, M., Najmul Islam, A., e Mäntymäki, M. (2022). How to explain ai systems to end users: a systematic literature review and research agenda. *Internet Research*, 32(7):1–31.

Morley, J., Floridi, L., Kinsey, L., e Elhalal, A. (2020). From what to how: an initial review of publicly available ai ethics tools, methods and research to translate principles into practices. *Science and engineering ethics*, 26(4):2141–2168.

Nissenbaum, H. (1996). Accountability in a computerized society. *Science and engineering ethics*, 2:25–42.

Nissenbaum, H. (2001). How computer systems embody values. *Computer*, 34(3):120–119.

Nunes, J. L., Barbosa, G. D., de Souza, C. S., e Barbosa, S. D. (2024). Using model cards for ethical reflection on machine learning models: an interview-based study. *Journal on Interactive Systems*, 15(1):1–19.

Ozkaya, I. (2019). Ethics is a software design concern. *IEEE Software*, 36(3):4–8.

Ozmen Garibay, O., Winslow, B., Andolina, S., Antona, M., Bodenschatz, A., Coursaris, C., Falco, G., Fiore, S. M., Garibay, I., Grieman, K., et al. (2023). Six human-centered artificial intelligence grand challenges. *International Journal of Human–Computer Interaction*, 39(3):391–437.

Pasricha, S. (2022). Ai ethics in smart healthcare. *IEEE Consumer Electronics Magazine*, 12(4):12–20.

Pereira, F. H., Prates, R. O., Maciel, C., e Pereira, V. (2016). Análise de interação antecipada e aspectos volitivos em sistemas de comunicação digital póstuma. In *XV Simpósio Brasileiro sobre Fatores Humanos em Sistemas Computacionais*.

Pereira, R., Darin, T., e Silveira, M. S. (2024). GranDIHC-BR: Grand Research Challenges in Human-Computer Interaction in Brazil for 2025-2035. In *Proceedings of the XXIII Brazilian Symposium on Human Factors in Computing Systems (IHC '24)*, New York, NY, USA. Association for Computing Machinery.

Prates, R. O., Rosson, M. B., e de Souza, C. S. (2015). Making decisions about digital legacy with google's inactive account manager. In *Human-Computer Interaction–INTERACT 2015: 15th IFIP TC 13 International Conference, Bamberg, Germany, September 14-18, 2015, Proceedings, Part I 15*, pages 201–209. Springer.

Prem, E. (2023). From ethical ai frameworks to tools: a review of approaches. *AI and Ethics*, 3(3):699–716.

Rodrigues, K. R. d. H., Carvalho, L. P., Pimentel, M. d. G. C., e Freire, A. P. (2024). GranDIHC-BR 2025-2035 - GC2: Ethics and Responsibility: Principles, Regulations, and Societal Implications of Human Participation in HCI Research. In *Proceedings of the XXIII Brazilian Symposium on Human Factors in Computing Systems (IHC '24)*, New York, NY, USA. ACM.

Shneiderman, B. (2020). Bridging the gap between ethics and practice: guidelines for reliable, safe, and trustworthy human-centered ai systems. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 10(4):1–31.

Umbrello, S. e van de Poel, I. (2021). Mapping value sensitive design onto ai for social good principles. *AI and Ethics*, 1(3):283–296.

UNESCO, C. Recommendation on the ethics of artificial intelligence.

Vainio-Pekka, H., Agbese, M. O. O., Jantunen, M., Vakkuri, V., Mikkonen, T., Rousi, R., e Abrahamsson, P. (2023). The role of explainable ai in the research field of ai ethics. volume 13, pages 1–39.

Valério, F. A., Guimarães, T. G., Prates, R. O., e Candello, H. (2017). Here's what i can do: Chatbots' strategies to convey their features to users. In *Proceedings of the xvi brazilian symposium on human factors in computing systems*, pages 1–10.

van Berkel, N., Tag, B., Goncalves, J., e Hosio, S. (2022). Human-centered artificial intelligence: a contextual morality perspective. In *Behaviour & Information Technology*, volume 41, pages 502–518.

Vilone, G. e Longo, L. (2021). Classification of explainable artificial intelligence methods through their output formats. *Machine Learning and Knowledge Extraction*, 3(3):615–661.

Xu, W. (2019). Toward human-centered ai: a perspective from human-computer interaction. *interactions*, 26(4):42–46.

Zack, T., Lehman, E., Suzgun, M., Rodriguez, J. A., Celi, L. A., Gichoya, J., Jurafsky, D., Szolovits, P., Bates, D. W., Abdulnour, R.-E. E., et al. (2024). Assessing the potential of gpt-4 to perpetuate racial and gender biases in health care: a model evaluation study. *The Lancet Digital Health*, 6(1):e12–e22.