

TERMS OF  
(DIS)SERVICE

**Comparing  
misinformation policies  
in text-generative  
AI chatbot**

FEBRUARY 2025

Author: Maria Giovanna Sessa (EU DisinfoLab)

Reviewers: Amaury Lesplingart (Check First), Joe McNamee (EU DisinfoLab)

# INTRODUCTION

- Large language models (LLMs) are rapidly proliferating. Like any technological tool, they can be harnessed for legitimate purposes but also misused or exploited by malicious actors. As these models become more integrated into everyday applications, concerns about their role in spreading misinformation continue to grow, calling for solid policies to prevent this threat.
- This factsheet collects and analyses the misinformation-related policies of 11 leading chatbots, selected based on [NewsGuard](#)'s selection. Our focus is on text-generative AI, given its widespread use across various domains – including content creation, translation, and summarisation – as well as its role in assisting users by answering questions.
- For each LLM, we examine key policy elements, including explicit references to misinformation and related prohibited activities – such as scams or impersonation. Additionally, we outline content moderation practices, user reporting mechanisms, and the consequences of violating the platform's Terms of Service (ToS).
- A note on methodology and limitations: the information provided here is based on publicly available sources, which vary in clarity, accessibility, and format. While we have made every effort to compile a thorough and accurate guide, gaps or omissions may exist due to the availability of information at the time of our writing. If any inaccuracies are identified, we welcome feedback and will gladly make corrections.

TABLE 1

General policy references, legal compliance, and mentions of misinformation

Platform	Policy reference	Mention of misinformation
Chat GPT (OpenAI)	<ul style="list-style-type: none"> <li>• <a href="#">Usage policies</a> (updated on 29/01/2025).</li> <li>• Law compliance statement.</li> </ul>	<p>Yes</p> <p><i>Universal Policies:</i></p> <ul style="list-style-type: none"> <li>• Prohibition to harm others (including scams and misleading activities).</li> </ul> <p><i>Building with ChatGPT:</i></p> <ul style="list-style-type: none"> <li>• Prohibition to misinform, misrepresent, or mislead others (including disinformation and impersonation).</li> </ul>
Claude (Anthropic)	<ul style="list-style-type: none"> <li>• <a href="#">Consumer Terms of Service</a> (updated on 13/05/24).</li> <li>• Law compliance statement.</li> </ul>	<p>Yes</p> <p><i>Acceptable Use Policy:</i></p> <ul style="list-style-type: none"> <li>• Prohibition to spread misinformation.</li> <li>• Prohibition to create political campaigns or interfere in elections.</li> <li>• Prohibition to engage in fraudulent, abusive, predatory practices.</li> <li>• Prohibition to abuse the platform (e.g., coordinating malicious activities).</li> </ul>
Copilot (Microsoft)	<ul style="list-style-type: none"> <li>• <a href="#">Copilot AI Experiences Terms</a> (updated on 1/10/2024).</li> <li>• Law compliance statement.</li> <li>• Code of Conduct signatory.</li> </ul>	<p>Yes</p> <p><i>Code of Conduct:</i></p> <ul style="list-style-type: none"> <li>• Prohibition to engage in fraudulent, false, or misleading activities (including disinformation, and impersonation).</li> </ul>
Deepseek-R1 (DeepSeek)	<ul style="list-style-type: none"> <li>• <a href="#">DeepSeek Terms of Use</a> (updated on 20/01/2025).</li> <li>• Law compliance statement.</li> </ul>	<p>No</p> <p><i>Terms of Use:</i></p> <ul style="list-style-type: none"> <li>• No explicit reference to misinformation.</li> <li>• Prohibition to harm and impersonate others.</li> </ul>
Gemini (Google)	<ul style="list-style-type: none"> <li>• <a href="#">Terms of Service</a> (Google) (updated on 22/05/2024).</li> <li>• <a href="#">Generative AI-prohibited use policy</a> (updated on 17/12/2024).</li> <li>• <a href="#">Gemini API Additional Terms of Service</a> (updated on 5/02/2025).</li> <li>• Law compliance statement.</li> <li>• Code of Conduct signatory.</li> </ul>	<p>Yes</p> <p><i>Generative AI-prohibited Use Policy:</i></p> <ul style="list-style-type: none"> <li>• Prohibition to misinform, misrepresent, or conduct misleading activities (including scams and impersonation).</li> </ul>

		<ul style="list-style-type: none"> <li>• Exceptions<sup>1</sup> can be made “based on educational, documentary, scientific, or artistic considerations” or where benefits outweigh harms.</li> </ul>
Grok-2 (xAI)	<ul style="list-style-type: none"> <li>• <a href="#">Terms of Service – Consumer</a> (updated on 2/01/2025).</li> <li>• <a href="#">Terms of Service – Enterprise</a> (updated on 23/12/2024).</li> <li>• Law compliance statement.</li> </ul>	<p>No</p> <p><i>xAI Acceptable Use Policy:</i></p> <ul style="list-style-type: none"> <li>• No explicit reference to misinformation.</li> <li>• Prohibition to scam.</li> <li>• Prohibition to mislead.</li> </ul>
Le Chat (Mistral AI)	<ul style="list-style-type: none"> <li>• <a href="#">Legal terms and conditions – Terms of Service</a> (updated on 6/02/2025).</li> <li>• Law compliance statement.</li> </ul>	<p>Yes</p> <p><i>Use Policy:</i></p> <ul style="list-style-type: none"> <li>• Prohibition to engage in fraudulent activities or scams (including impersonation).</li> <li>• Prohibition to misinform.</li> </ul>
Meta AI (Meta)	<ul style="list-style-type: none"> <li>• <a href="#">Meta AIs Terms of Service (EU)</a><sup>2</sup> (updated on 17/01/2025).</li> <li>• <a href="#">Meta AIs Terms of Service</a> (updated on 17/01/2025).</li> <li>• Law compliance statement.</li> <li>• Code of Conduct signatory.</li> </ul>	<p>Yes</p> <p><i>Meta AIs Terms of Service (EU):</i></p> <ul style="list-style-type: none"> <li>• Prohibition to deceive or mislead others (including scams, misinformation, or disinformation).</li> </ul>
Perplexity API (Perplexity)	<ul style="list-style-type: none"> <li>• <a href="#">Perplexity API Terms of Service</a> (updated on 17/02/2024).</li> <li>• Law compliance statement.</li> </ul>	<p>No</p> <p><i>Use of platform:</i></p> <ul style="list-style-type: none"> <li>• No explicit reference to misinformation.</li> <li>• Prohibition to impersonate others or provide false information.</li> </ul>
Pi (Inflection)	<ul style="list-style-type: none"> <li>• <a href="#">Privacy Policy</a> (updated on 19/09/2023).</li> <li>• Law compliance statement.</li> </ul>	<p>Yes</p> <p><i>Acceptable Use:</i></p> <ul style="list-style-type: none"> <li>• Prohibition to engage in harmful uses (including misinformation).</li> <li>• Prohibition to impersonate others during the sign-up process</li> </ul>
You.com	<ul style="list-style-type: none"> <li>• <a href="#">Terms of Service</a> (updated on 27/08/2024).</li> <li>• Law compliance statement.</li> </ul>	<p>Yes</p> <p><i>Acceptable Use Policy:</i></p> <ul style="list-style-type: none"> <li>• Prohibition to misinform, deceive, or mislead others (including disinformation, misinformation, and impersonation).</li> </ul>

<sup>1</sup> It is unclear who makes these decisions concerning exceptions and when.

<sup>2</sup> We find it contradictory that Meta has EU-specific Meta AI ToS despite being unavailable in the EU.

TABLE 2

## ToS violations: detection, moderation, consequences, reporting, and appeal mechanisms

Platform	ToS violation detection and content moderation	ToS violation reporting	Consequences for ToS violation and appeal mechanism <sup>3</sup>
Chat GPT (OpenAI)	<ul style="list-style-type: none"> <li>• <a href="#">Combination</a> of automated and human detection, and user reports.</li> </ul>	<ul style="list-style-type: none"> <li>• Product interface reporting using the “thumbs down”.</li> <li>• <a href="#">Content reporting form</a>.</li> </ul>	<ul style="list-style-type: none"> <li>• Actions against content or account (e.g., warnings, restrictions, ineligibility for GPT Store or monetisation), suspension, or termination.</li> <li>• <a href="#">Account ban appeal</a>.</li> </ul>
Claude (Anthropic)	<ul style="list-style-type: none"> <li>• <a href="#">Content moderation</a> by Trust and Safety team (process unspecified).</li> <li>• Adjustable <a href="#">safety filters</a> on prompts.</li> <li>• User’s responsibility to verify output accuracy independently.</li> </ul>	<ul style="list-style-type: none"> <li>• Product interface <a href="#">reporting</a> to Trust &amp; Safety team using the “thumbs down” button.</li> <li>• Reporting via email (usersafety@anthropic.com).</li> </ul>	<ul style="list-style-type: none"> <li>• Account throttling, suspension, or termination.</li> <li>• <a href="#">Account ban appeal</a>.</li> </ul>
Copilot (Microsoft)	<ul style="list-style-type: none"> <li>• Content restriction or removal if in violation of the <a href="#">Code of Conduct</a> (process unspecified).</li> <li>• User’s responsibility to verify output accuracy independently.</li> </ul>	<ul style="list-style-type: none"> <li>• Product interface reporting using the Feedback button.</li> <li>• <a href="#">Content reporting form</a>.</li> </ul>	<ul style="list-style-type: none"> <li>• Account suspension, service limitations, or content removal, deletion, or restriction.</li> <li>• <a href="#">Account ban appeal</a>.</li> </ul>
Deepseek-R1 (DeepSeek)	<ul style="list-style-type: none"> <li>• Content removal (process unspecified).</li> <li>• User’s responsibility to verify output accuracy independently.</li> </ul>	<ul style="list-style-type: none"> <li>• Product interface reporting using the “Contact Us” button.</li> <li>• Reporting via email(chat: service@deepseek.com; platform: api-service@deepseek.com) or contact address (5th Floor, North Building, Block C, Rongke Information Center, No.2 South Science Academy Road, Haidian District, Beijing, China).</li> </ul>	<ul style="list-style-type: none"> <li>• Warnings (and deadline for correction), account function restriction, suspension, termination, prohibition to re-register, and content removal.</li> <li>• <a href="#">Account ban appeal</a>.</li> </ul>
Gemini (Google)	<ul style="list-style-type: none"> <li>• <a href="#">Combination</a> of automated and human detection by</li> </ul>	<ul style="list-style-type: none"> <li>• Product interface <a href="#">reporting</a> using the “Bad response” button or</li> </ul>	<ul style="list-style-type: none"> <li>• Content removal.</li> </ul>

<sup>3</sup> The table reports the language used by each platform. Due to the lack of definitions provided, we cannot determine whether some terms – such as ‘termination’ and ‘deletion’ are synonymous.

	<p>Google’s Trust and Safety Team.</p> <ul style="list-style-type: none"> <li>Adjustable <a href="#">safety filters</a> on prompts.</li> </ul>	<p>selecting “Help” and then “Report a problem”.</p> <ul style="list-style-type: none"> <li><a href="#">Content reporting form</a>.</li> </ul>	<ul style="list-style-type: none"> <li>Google account suspension or termination.</li> <li>Google <a href="#">account ban appeal</a>.<sup>4</sup></li> </ul>
Grok-2 (xAI)	<ul style="list-style-type: none"> <li>Content removal (process unspecified).</li> <li>User’s responsibility to verify output accuracy independently.</li> </ul>	<ul style="list-style-type: none"> <li>No reporting mechanism.</li> </ul>	<ul style="list-style-type: none"> <li>Content deletion or disabling.</li> <li>Service access suspension or termination, account deletion.</li> <li>Account ban appeal via email (support@x.ai).</li> </ul>
Le Chat (Mistral AI)	<ul style="list-style-type: none"> <li>Automated <a href="#">monitoring</a>.</li> <li>User’s responsibility to verify output accuracy independently.</li> </ul>	<ul style="list-style-type: none"> <li>Product interface reporting using the “thumbs down” feature.</li> <li>Reporting via <a href="#">Help Centre</a> or email (legal@mistral.ai).</li> </ul>	<ul style="list-style-type: none"> <li>Account termination or suspension.</li> <li>Account ban appeal via the <a href="#">Help Centre</a> or email (support@mistral.ai).</li> </ul>
Meta AI (Meta)	<ul style="list-style-type: none"> <li><a href="#">Combination</a> of automated and human detection.</li> <li>User’s responsibility to verify output accuracy independently.</li> </ul>	<ul style="list-style-type: none"> <li>Product interface <a href="#">reporting</a> using the “thumbs down” button or selecting “Report a Problem”.</li> </ul>	<ul style="list-style-type: none"> <li>Account suspension or permanent disabling.</li> <li>Facebook <a href="#">account ban appeal</a>.<sup>5</sup></li> </ul>
Perplexity API (Perplexity)	<ul style="list-style-type: none"> <li>Content removal (process unspecified).</li> </ul>	<ul style="list-style-type: none"> <li>Product interface <a href="#">reporting</a> using the “flag icon”.</li> <li>Reporting via email (support@perplexity.ai) or submission form.</li> </ul>	<ul style="list-style-type: none"> <li>Account suspension or termination.</li> <li>No specific account ban appeal process but inquiries via email (support@perplexity.ai).</li> </ul>
Pi (Inflection)	<ul style="list-style-type: none"> <li>Content moderation (process unspecified) and removal (copyright infringement).</li> <li>User’s responsibility to verify output accuracy independently.</li> </ul>	<ul style="list-style-type: none"> <li>Product interface “Report” button (called the “Feedback”).</li> </ul>	<ul style="list-style-type: none"> <li>Account suspension or termination.</li> <li>Account ban appeal via email (support@pi.ai).</li> </ul>
You.com (You.com)	<ul style="list-style-type: none"> <li>Combination of automated and human detection.<sup>6</sup></li> </ul>	<ul style="list-style-type: none"> <li>Product interface reporting using the “thumbs down” feature or selecting “More” and then “Send Feedback”.</li> </ul>	<ul style="list-style-type: none"> <li>Service access suspension or termination, account termination.</li> <li>No specific account ban appeal process but inquiries via <a href="#">“Contact Us”</a>.</li> </ul>

<sup>4</sup> Since Gemini is linked to a user’s Google account, losing access to the Google account would mean losing access to Gemini too. However, it is unclear whether Google can suspend access to Gemini specifically while keeping the rest of the Google account functional, as it does with some services like YouTube and AdSense.

<sup>5</sup> As of now, it is unclear whether a Meta AI account is separate from a user’s primary Meta account (e.g., Facebook). Due to this ambiguity and the lack of relevant documentation, we referenced the Facebook account suspension process.

<sup>6</sup> This answer was provided by the You.com chatbot. However, we could not find evidence to support it, and when asked to reference its statement, the chatbot responded: “Content moderation policies, especially with AI-powered platforms like You.com, are often evolving and may not always be fully public.”

TABLE 3

## A comparative overview of text-generative LLMs' misinformation policies

	Law compliance statement	Code of Conduct signatory	Mention of misinformation	ToS violation detection and content moderation	Mention of user fact-checking responsibility	ToS violation reporting mechanism	Suspension and termination as ToS violation consequences	ToS violation appeal mechanism
Chat GPT (OpenAI)	Yes	No	Yes	Yes (automated + human)	No	Yes	Yes	Yes
Claude (Anthropic)	Yes	No	Yes	Yes (unspecified)	Yes	Yes	Yes	Yes
Copilot (Microsoft)	Yes	Yes	Yes	Yes (unspecified)	Yes	Yes	Yes	Yes
Deepseek-R1 (DeepSeek)	Yes	No	No	Yes (unspecified)	Yes	Yes	Yes	Yes
Gemini (Google)	Yes	Yes	Yes	Yes (automated + human)	No	Yes	Yes	Yes
Grok-2 (xAI)	Yes	No	No	Yes (unspecified)	Yes	No	Yes	Yes
Le Chat (Mistral AI)	Yes	No	Yes	Yes (automated)	Yes	Yes	Yes	Yes
Meta AI (Meta)	Yes	Yes	Yes	Yes (automated + human)	Yes	Yes	Yes	Unclear
Perplexity API (Perplexity)	Yes	No	No	Yes (unspecified)	No	Yes	Yes	No
Pi (Inflection)	Yes	No	Yes	Yes (unspecified)	Yes	Yes	Yes	Yes
You.com	Yes	No	Yes	Yes (automated + human)	No	Yes	Yes	No

## CONCLUDING REMARKS

- In the 11 chatbots considered, references to misinformation vary (and sometimes are missing). Even when the term is mentioned, it is rarely defined. This is not uncommon: what some platforms call ‘deletion’, others call ‘termination’, though potential nuances remain unclear. Additionally, it is unspecified who enforces practices, how long suspensions last, and what distinguished warnings, suspensions, and terminations. On another note, more stringent prohibitions generally focus on hate speech, highlighting the illegal nature of the issue.
- Overall, ToS violation detection and content moderation processes are unclear. Some platforms mention dedicated moderation teams and processes – typically a mix of automated and human reviews – but the details remain vague. There is no transparency regarding the resources invested in content moderation, including the number of personnel, their locations and language coverage, and the timeline for action.
- Given our focus on misinformation, we concentrated on content reporting. Nevertheless, it is evident that platforms prioritise reporting breaches of and vulnerabilities in their architecture over content-related concerns.
- Law compliance is usually referenced in broad terms as adherence to “all applicable legislation”. Platforms primarily emphasise privacy and data protection, the prohibition of illegal activities (with a particular focus on child protection), restrictions on reverse-engineering platform functions, and using the services for biometric identification. Notably, there is a perception that users bear greater responsibility for providing accurate data and avoiding platform abuse than the platforms themselves do in terms of accountability. This impression arises because the ToS are detailed regarding user obligations but more superficial about platform responsibilities. Uncertainty remains regarding platform incentives and their actual capacity to enforce these rules.
- Most platforms specify that fact-checking is the user’s responsibility as if trying to exempt themselves from liability related to providing false and inaccurate answers. However, it is worth noting that the general provisions of the EU’s Digital Services Act (DSA) require online platforms to remove illegal content expeditiously once they have actual knowledge of its illegality.
- If current platform approaches to misinformation already seem inadequate, there are valid concerns that the situation will deteriorate further. The cross-platform trend of dismantling safety policies and moderation teams reinforces these fears. Examples pointing in this direction include OpenAI’s decision to remove some content warnings on ChatGPT, and xAI’s introduction of “Unhinged Mode” in Grok-2, a beta feature described as intentionally “objectionable, inappropriate, and offensive, much like an amateur stand-up comic who is still learning the craft.”



EU DISINFO LAB

[www.disinfo.eu](http://www.disinfo.eu)

