

# Governing hate speech detection on online platforms: A human rights approach

Eva Nave<sup>1</sup> and Lottie Lane<sup>2</sup>

## Abstract

The Internet is a global forum largely governed by private actors driven by profit concerns, often disregarding the human rights of historically marginalised communities. Increased attention is being paid to the corporate human rights due diligence (HRDD) responsibilities applicable to online platforms countering illegal online content, such as hate speech. At the European Union (EU) level, cross-sector initiatives regulate the rights of marginalised groups and establish HRDD responsibilities for online platforms to expeditiously identify, prevent, mitigate, remedy and remove online hate speech. These initiatives include the Digital Services Act, the Audiovisual Media Services Directive, the Directive on Corporate Sustainability Due Diligence, the Artificial Intelligence Act and the Code of conduct on countering illegal hate speech online. Nevertheless, the HRDD framework applicable to online hate speech has focused mostly on the platforms' responsibilities throughout the course of their operations - guidance regarding HRDD requirements concerning the regulation of hate speech in the platforms' Terms of Service (ToS) is missing. This paper<sup>3</sup> critically employs a conceptualisation of criminal hate speech as explained in the Council of Europe Committee of Ministers' Recommendation CM/Rec(2022)16, Paragraph 11, to develop specific HRDD responsibilities. We argue that online platforms should, as part of emerging preventive HRDD responsibilities within Europe, respect the rights of historically oppressed communities by aligning their ToS with the conceptualisation of criminal hate speech in European human rights standards.

## 1 Introduction

This paper addresses a vacuum in the legal framework by clarifying corporate human rights responsibilities in Europe to counter the most serious forms of online hate speech. In particular, we examine (emerging) standards on human rights due diligence (HRDD), artificial intelligence and online content moderation at the international and European level. We claim that there is a legal standard emanating from the HRDD framework in the European context prescribing the responsibility for online platforms - particularly for very large online platforms, video-sharing platforms and for platforms under the scope of the Directive on corporate sustainability due diligence<sup>4</sup> - to align their terms of service with the conceptualization of the criminal hate speech in Paragraph 11 of the Council of Europe Committee of Ministers Recommendation CM/Rec(2022)16.<sup>5</sup> Based on this claim, we provide recommendations to law- and policy-makers.

Around two thirds of the world's population are active Internet users.<sup>6</sup> While the Internet enables individuals to access information and exercise their freedom of expression, it also enables the proliferation of online hate speech. 'Online hate speech' broadly refers to discriminatory expressions shared on digital environments targeting historically or systematically marginalized<sup>7</sup> people. Recommendation CM/Rec(2022)16 reiterates that hateful expressions represent a violation of human rights. When unaddressed, these can hinder peace and development by denying the values of pluralism, tolerance and broadmindedness essential in a democratic society.

The rise of online hate speech results from specific features of the Internet. First, unlike in traditional media, most con-

---

<sup>1</sup> PhD candidate, Center for Law and Digital Technologies, Leiden University, <https://orcid.org/0000-0001-5040-1601>, [e.v.r.nave@law.leidenuniv](mailto:e.v.r.nave@law.leidenuniv)

<sup>2</sup> Assistant Professor of Public International Law, University of Groningen, <https://orcid.org/0000-0002-2986-3075>, [c.l.lane@rug.nl](mailto:c.l.lane@rug.nl)

<sup>3</sup> The completed version of this paper is published in the Computer Law and Security Review, Special Issue on Vulnerability, Marginalization and Data Protection Law, available at: <https://doi.org/10.1016/j.clsr.2023.105884>

<sup>4</sup> European Commission (2022) Proposal for a Directive of the European Parliament and of the Council on Corporate Sustainability Due Diligence and amending Directive (EU) 2019/1937.

<sup>5</sup> Council of Europe Committee of Ministers, Recommendation CM/Rec(2022)16 of the Committee of Ministers to member States on combating hate speech (CM/Rec(2022)16).

<sup>6</sup> Number of internet and social media users worldwide as of January 2023 (2023) <<https://www.statista.com/statistics/617136/digital-population-worldwide/>> accessed 26 April 2024.

<sup>7</sup> This paper uses 'oppression' and 'marginalisation' interchangeably.

68 tent published on the Internet can be quickly shared with lit-  
69 tle to no monitoring, made available to large audiences, pub-  
70 lished under anonymity, and easily manipulated in ways that  
71 intensify hate (e.g. hate profiles, memes and deep fakes).  
72 Second, online content is hosted by businesses primarily  
73 driven by profit goals, often at the expense of human rights.  
74 The potentially negative impact of AI-driven content moder-  
75 ation by online platforms is under increasing scrutiny. For  
76 example, Meta Platforms, Inc. (formerly named Facebook,  
77 Inc.) faces legal action for alleged negligence in facilitating  
78 the genocide of Rohingya Muslims in Myanmar after its al-  
79 gorithm failed to remove hateful posts and amplified hate  
80 speech.<sup>8</sup> Similarly, whistle-blower Frances Haugen alerted  
81 that Facebook neglected reports of accounts and hate speech  
82 content towards Muslims in India, potentially leading to of-  
83 fline violence.<sup>9</sup> There are reportedly other situations of hu-  
84 man rights abuses by different platforms.<sup>10</sup>

## 86 2 Responses by scholars, legislators and policy- 87 makers

88 Legal scholars have emphasized the growing impact of so-  
89 cial media platforms in the application of regulatory frame-  
90 works for freedom of expression and democratic processes  
91 and the subsequent need to expand the legal scholarship fo-

92 cusing on the regulation of online platforms.<sup>11</sup> In this con-  
93 text, it is relevant to consider that most online platforms are  
94 based in the USA and thus typically bound by the USA  
95 framework on freedom of expression, corporate human  
96 rights due diligence and intermediary liability. Conversely,  
97 to the extent that online platforms operate in European Union  
98 (EU) territory, they must also abide by the regional human  
99 rights frameworks in Europe, which differ significantly from  
100 those applicable in the USA. The reconciliation of different  
101 regional standards has been challenging, not only for online  
102 platforms but also for judicial bodies in enforcing their deci-  
103 sions.<sup>12</sup>

105 Legislators and policy-makers at the international, regional  
106 and national level have made efforts to prevent and address  
107 the negative impact of business on human rights, including  
108 through HRDD and through liability regimes.<sup>13</sup> The HRDD  
109 regime includes the seminal United Nations Guiding Princi-  
110 ples on Business and Human Rights (UNGPs), which are ar-  
111 guably the most authoritative international expression of the  
112 corporate responsibility to respect human rights through

---

<sup>8</sup> Al Jazeera, 'Rohingya sue Facebook for \$150bn for fuelling Myanmar hate speech' (7 December 2021) <<https://www.aljazeera.com/news/2021/12/7/rohingya-sue-facebook-for-150bn-for-fuelling-myanmar-hate-speech>> accessed 26 April 2024.

<sup>9</sup> Al Jazeera, 'Facebook failing to check hate speech, fake news in India: Report' (25 October 2021) <<https://www.aljazeera.com/news/2021/10/25/facebook-india-hate-speech-misinformation-muslims-social-media>> accessed 26 April 2024.

<sup>10</sup> Shaun Harper, 'Hate Speech Rises On Twitter After Elon Musk Takes Over, Researchers Find' (*Forbes*, 31 October 2022) <<https://www.forbes.com/sites/shaunharper/2022/10/31/elon-musk-twitter-takeover-leads-to-n-word-and-hate-speech-increase-lebron-james-calls-for-action/?sh=f28a381dd99a>> accessed 6 April 2023; Hadi Al Khatib and Dia Kayyali, 'YouTube Is Erasing History' (*The New York Times*, 23 October 2019) <<https://www.nytimes.com/2019/10/23/opinion/syria-youtube-content-moderation.html>> accessed 26 April 2024.

<sup>11</sup> E.g. Kate Klonick, 'The new governors: The people, rules, and processes governing online speech' (2017) *Harv. L. Rev.*, 131, 1598; Tarlach McGonagle, 'Free Expression and Internet Intermediaries: The Changing Geometry of European Regulation' (2020) *Oxford Handbooks in Law* (pp. 467–485), 10; Tarlach McGonagle, 'The Council of Europe and Internet Intermediaries: A Case Study of Tentative Posturing', 232, in Rikke Frank Jørgensen (eds), 'Human Rights in the Age of Platforms' (2019) Cambridge, MA: The MIT Press <<https://doi.org/10.7551/mitpress/11304.001.0001>> accessed 26 April 2024; Judit

Bayer, Bernd Holznel, Päivi Korpisaari (ex. Tiilikka), Lorna Woods, Volume 1' (2021) Baden-Baden: Nomos Verlagsgesellschaft mbH & Co. KG., 30, <<https://doi.org/10.5771/9783748929789>> accessed 26 April 2024; Martin Moore and Tambini Damian (eds), 'Regulating Big Tech: Policy Responses to Digital Dominance' (2021), <<https://doi.org/10.1093/oso/9780197616093.001.0001>> accessed 26 April 2024.

<sup>12</sup> E.g. *Ligue contre le racisme et l'antisémitisme et Union des étudiants juifs de France c. Yahoo! Inc. et Société Yahoo! France (LICRA v. Yahoo!)*; and *Yahoo Inc. v LICRA*; European Court of Justice, Opinion of Advocate General Szpunar delivered on 8 June 2023 (1) Case C-376/22 clarifies that Union law prescribes the possibility for Member States to restrict the freedom to provide information society services to 'fight against any incitement to hatred on grounds of race, sex, religion or nationality, and violations of human dignity concerning individual persons'.

<sup>13</sup> Council of Europe, Committee on Artificial Intelligence, Consolidated Working Draft of the Framework Convention on Artificial Intelligence, Human Rights, Democracy and the Rule of Law, CAI(2023)18 <<https://rm.coe.int/cai-2023-18-consolidated-working-draft-framework-convention/1680abde66>> accessed 26 April 2024. Though State-centric and though not including standards directly applicable to companies, this Framework Convention provides key guidance for States regarding human rights centred approached to the governance and regulation of artificial intelligence.

113 HRDD.<sup>14</sup> At the EU level, the Directive on corporate sustain-  
114 ability due diligence (CSDDD) was just recently adopted.<sup>15</sup>  
115 Businesses - including online platforms - falling under the  
116 scope of the proposal should identify, prevent, mitigate and  
117 bring an end to negative impacts on human rights. Further-  
118 more, the EU Artificial Intelligence Act (AI Act) emphasizes  
119 the need for protection of human rights in the digital envi-  
120 ronment.<sup>16</sup>

121  
122 Concerning HRDD and moderation of harmful content  
123 online, in November 2022 the Regulation for a Digital Ser-  
124 vices Act entered into force.<sup>17</sup> The Digital Services Act adds  
125 to the EU Audiovisual Media Services Di-  
126 rective (AVMSD)<sup>18</sup> and enhances cross-sector due diligence  
127 responsibilities for digital services to remove illegal content  
128 online. This includes hate speech.<sup>19</sup> The due diligence frame-  
129 work in the Digital Services Act aligns with  
130 CM/Rec(2022)16 and builds on the Code of conduct on  
131 countering illegal hate speech online whereby IT companies  
132 commit to expeditiously review and remove hate speech and  
133 to promote transparency towards users.<sup>20</sup>

### 135 3 Gaps in law and policy

136 Despite these advancements, the HRDD framework applica-  
137 ble to online hate speech has focused mostly on explaining  
138 the responsibilities of companies throughout their opera-  
139 tions. Guidance regarding HRDD requirements for the regu-  
140 lation of hate speech in the terms of service is missing. A key  
141 aspect remains unaddressed: how online platforms should  
142 define hate speech and how this should be communicated to  
143 their users. More specifically, is there a legal standard ema-  
144 nating from the European HRDD framework prescribing the  
145 responsibility for online platforms<sup>21</sup> to align their terms of  
146 service, as a minimum legal standard, with the conceptual-

147 ization of the criminal hate speech as explained in the Euro-  
148 pean human rights standards, in particular with the Recom-  
149 mendation CM/Rec(2022)16?  
150

### 151 4 Scope of research

152 To answer this research question, we employ doctrinal re-  
153 search to differentiate between hate speech that is criminally  
154 actionable and hate speech that should be prohibited under  
155 civil or administrative law (Section 2). We focus on hate  
156 speech that is criminally actionable. The limitation of the re-  
157 quirement to harmonize and reflect the conceptualization of  
158 *criminal* hate speech is justified by a growing human rights  
159 understanding of criminal hate speech as reflected in  
160 CM/Rec/(2022)16, Para. 11, from which specific HRDD re-  
161 sponsibilities can be developed. Building upon critical race,  
162 black feminist and intersectionality theories, this paper  
163 claims that platforms should explicitly conceptualize hate  
164 speech as discriminatory communications that target an  
165 open-ended list of historically or systematically oppressed  
166 people or group of people. This conceptualization should  
167 also consider the intersectionality of systems of oppression  
168 as an aggravating harm resulting from hate speech. It is  
169 worth remembering that the European Commission proposed  
170 to add hate speech to the list of EU crimes which, if and when  
171 this proposal materializes, will strengthen the need for a  
172 standardized conceptualization of criminal hate speech in  
173 online platforms' terms of service. This legal avenue sup-  
174 ports compliance with the transparency and clarity required  
175 on Terms and Conditions (Article 14 Digital Services Act)  
176 generally imposed on all providers of intermediary services.  
177

### 178 5 Key findings

179 In Sections 3 and 4, we investigate the HRDD regime.<sup>22</sup> We  
180 examine the HRDD framework applicable to AI businesses

---

<sup>14</sup> UN Human Rights Council, 'Report of the Special Representative of the Secretary-General on the Issue of Human Rights and Transnational Corporations and Other Business Enterprises, John Ruggie' (2011) A/HRC/17/31. We use the term 'responsibility' to denote non-legally binding standards and 'obligation' when discussing binding standards.

<sup>15</sup> European Commission (2022) Proposal for a Directive of the European Parliament and of the Council on Corporate Sustainability Due Diligence and amending Directive (EU) 2019/1937.

<sup>16</sup> European Commission (2021) Proposal for a Regulation of the European Parliament and of the Council laying down harmonised rules on artificial intelligence (Artificial Intelligence Act) and amending certain Union legislative acts COM(2021) 206 final, Explanatory Memorandum, 1.1.

<sup>17</sup> European Union, Regulation of the European Parliament and of the Council on a Single Market For Digital Services (Digital Services Act) and amending Directive 2000/31/EC, Article 93.

<sup>18</sup> Directive 2010/13/EU of the European Parliament and of the Council of 10 March 2010 on the coordination of certain

provisions laid down by law, regulation or administrative action in Member States concerning the provision of audiovisual media services (Audiovisual Media Services Directive), OJ L 95.

<sup>19</sup> European Commission (2018) Commission Recommendation (EU) 2018/334 of 1 March 2018 on measures to effectively tackle illegal content online, OJ L 63/50.

<sup>20</sup> European Commission (2016) The Code of Conduct on countering illegal hate speech online.

<sup>21</sup> 'AI businesses' is sometimes used synonymously with 'IT companies' or by 'internet intermediaries' (or 'intermediaries'), depending on the legal instrument under analysis.

<sup>22</sup> As a regulatory approach distinct from that of HRDD – as seen in the separate chapters on each regime in the Digital Services Act –, the EU liability regime for internet service providers (ISP) falls outside the remit of this research. These regimes are nevertheless related in that liability may follow from non-compliance with HRDD responsibilities. For discussion of ISP liability regimes and recent case law, see e.g. Andrea Bertolini et al., 'Liability of Online Platform: Study for the European Parliament' (2021) European Parliamentary Research Service PE 656.318; Berrak Genç-Gelgeç,

181 (the UNGPs, initiatives of the Organization for Economic  
182 Cooperation and Development, the CSDDD and the AI Act)  
183 and preventive HRDD responsibilities in moderation of ille-  
184 gal content, such as criminal hate speech (with reference to  
185 the Digital Services Act, the AVMSD, the Code of Conduct<sup>23</sup>  
186 and CM/Rec(2022)16). On this basis, we argue that terms of  
187 service fulfil the role of the human rights ‘policy commit-  
188 ment’ standard found in the UNGPs and should include a  
189 clear explanation of the platforms’ commitment to human  
190 rights, including the prohibition of criminal hate speech. The  
191 drafting and updating of the terms of service should be a  
192 means for online platforms to respond to the systemic risk of  
193 online hate speech and terms of service should explicitly re-  
194 flect the HRDD responsibilities to prohibit, remove and re-  
195 port criminal hate speech to relevant public authorities.

196  
197 We suggest that to improve legal coherence in countering  
198 online hate speech in the European context, online platforms  
199 should follow CM/Rec(2022)16 Para.11’s conceptualization  
200 of criminal hate speech in their terms of service. Paragraph  
201 11 builds upon binding and non-binding international human  
202 rights standards, such as the International Convention on the  
203 Elimination of All Forms of Racial Discrimination,<sup>24</sup> the  
204 Convention on the Prevention and Punishment of the Crime  
205 of Genocide,<sup>25</sup> International Covenant on Civil and Political  
206 Rights,<sup>26</sup> Article 20(2), etc. As a result, some of these ele-  
207 ments can be claimed to broadly represent international hu-  
208 man rights standards and could therefore be extrapolated to  
209 international preventive HRDD responsibilities to counter  
210 criminal hate speech.

## 212 6 Assessment of practice

213 Section 5 presents an empirical qualitative analysis of three  
214 case studies: Facebook, Twitter and YouTube. We assess the  
215 compliance of the platforms’ terms of service with the Euro-  
216 pean Court of Human Rights jurisprudence on criminal hate  
217 speech, and with the conceptualization of criminal hate  
218 speech in CM/Rec(2022)16. The platforms were selected be-  
219 cause they: (1) fall under the scope of CSDDD; (2) are sig-  
220 natories to the EU Code of Conduct; and (2) qualify as very  
221 large online platforms (VLOPs) as defined in the Digital Ser-  
222 vices Act.<sup>27</sup> The three case studies demonstrate that although

223 Facebook, Twitter and YouTube have each to a certain de-  
224 gree adopted terms of service prohibiting hate speech, none  
225 of them currently conceptualizes hate speech in a way that is  
226 consistent with human rights standards. More specifically,  
227 none recognizes the difference between hate speech crimi-  
228 nally actionable and hate speech which may be prohibited  
229 under civil or administrative law. Moreover, none recognizes  
230 the specific HRDD responsibilities associated with counter-  
231 ing criminal hate speech. The three case studies reveal the  
232 lack of alignment of content moderation practices by online  
233 platforms with the HRDD responsibilities to identify, miti-  
234 gate, cease, remedy and inform about potentially adverse im-  
235 pacts on human rights.

## 237 7 Recommendations and conclusion

238 In summary, addressing law- and policy-makers, we recom-  
239 mend that the European Commission issues a best practice  
240 guideline (under Article 35(3) Digital Services Act and Arti-  
241 cle 13 CSDDD) suggesting that VLOPs, and particularly  
242 video-sharing platforms, should explicitly mention in their  
243 terms of service that they prohibit, remove and report to law  
244 enforcement authorities criminal hate speech in line with the  
245 conceptualization in Paragraph 11 CM/Rec(2022)16. Further  
246 to this and also by issuing a best practice guideline, we rec-  
247 ommend that the European Commission suggests that  
248 VLOPs, with a similar heightened focus on video-sharing  
249 platforms, adopt HRDD compliant content moderation pro-  
250 cesses which should likewise be explicitly mentioned in their  
251 terms of service.

252  
253 This paper has primarily addressed the first phase of HRDD  
254 processes, i.e. the adoption of a policy commitment as a pre-  
255 ventive HRDD responsibility. Further research is necessary  
256 to examine what could be required in relation to the remain-  
257 ing phases of HRDD, i.e. the tracking and communicating  
258 implementation and results as well as the provision of reme-  
259 dies when applicable. For example, what online platforms  
260 moderating content should do to identify and prevent the pro-  
261 motion of criminal hate speech, and how they could effec-  
262 tively respond to these risks, should be the subject of further  
263 study.

---

‘Regulating Digital Platforms: Will the DSA Correct Its Pre-  
decessor’s Deficiencies?’ (2022) 18 Croatian Yearbook of  
European Law and Policy 25; United States Supreme Court,  
Twitter v. Taamneh 598 US (2023).

<sup>23</sup> Some EU instruments use the problematic expression ‘ille-  
gal hate speech’, which could lead the reader to understand  
that there is legal hate speech. This is not the case. Hate  
speech is always illegal but it can be criminalised only in its  
most serious forms. For legal coherence purposes, this paper  
will refrain from using ‘illegal hate speech’ unless referring  
to the title of an instrument.

<sup>24</sup> UN General Assembly, International Convention on the  
Elimination of All Forms of Racial Discrimination, 21 De-  
cember 1965, United Nations, Treaty Series, vol. 660.

<sup>25</sup> UN General Assembly, Convention on the Prevention and  
Punishment of the Crime of Genocide, 9 December 1948,  
United Nations, Treaty Series, vol. 78.

<sup>26</sup> UN General Assembly, International Covenant on Civil and  
Political Rights, 16 December 1966, United Nations, Treaty  
Series, vol. 999.

<sup>27</sup> i.e. they have 45 million or more average monthly active  
recipients of their service in the Union: Digital Services Act,  
Recital 76.