

## Manuscript

### State Sovereignty, Disinformation and Speech Governance

*Julia Hörnle, Queen Mary University of London*

#### 1. Introduction

A powerful example of the ominous connection between disinformation and the dangers of unchecked social media companies is the recent decision by *Meta* to withdraw the system of fact-checkers from Facebook and Instagram in the US, and replace the fact-checking with community notes à la X, currying political favours<sup>1</sup> with the incoming Trump administration.<sup>2</sup> A political alliance between an incoming government and media organisations is not likely to improve free speech in respect of criticizing the new government, and it is undermining the role of media in holding those in power to account. The apparent reason given for this move is greater freedom of speech and removing subjectivity from fact-checking by news organisations.

At least from a European perspective, this argument is paradoxical: how can there be free speech if we do not ensure that political speech is based on factual truth? In other words, in a rational world, free speech without factual truth is useless. One important component of the right to free expression under a European conception is the right to receive information<sup>3</sup> and disinformation threaten this part of the equation. Of course, not all speech is grounded in facts, but where speech in news information sources<sup>4</sup> is grounded in facts, should we not ensure that the facts themselves are checked? Disinformation in Europe is regarded as a fundamental threat to freedom of expression as it endangers the information ecosystem. It is therefore seen as an existential threat to the liberal, democratic ordering of the state.<sup>5</sup>

The significance of fact-checking has to be put in the context of the ongoing ideology war in the US around free speech, disinformation and the increasing gulf between liberal Democrat and right-wing Republican world view and the corresponding perception of reality. This ideology war instrumentalizes free speech in the fight for power, starting with the so-called

---

<sup>1</sup>According to news reports, Meta has donated \$1 million to Trump's inauguration fund and replaced the President of Global Affairs Nick Clegg (a former British Liberal Democrat politician) with the Republican Joel Kaplan (a former White House Deputy Chief of Staff for Policy under George W Bush). Zuckerberg might hope that in return Trump is able to curb social media regulation in Brazil and Europe, for example in the context of trade negotiations.

<sup>2</sup> BBC News (7. January 2025) <https://www.bbc.co.uk/news/articles/cly74mpy8klo>; New York Times (7. January 2025) <https://www.nytimes.com/live/2025/01/07/business/meta-fact-checking>

<sup>3</sup> Art 10 ECHR

<sup>4</sup>According to Ofcom's research in the UK, just over half (52%) of UK adults use social media for obtaining news. Younger generations have turned away from traditional news sources- the 16-24 year old cohort use social media as a main source for news (82%). However, the detailed picture is more complicated, 70% of adults (but only 49% of 16-24 year olds) continue to rely on TV news, and traditional media have scored high for trust, accuracy and impartiality- see Ofcom, News Consumption in the UK (10. September 2024) <https://www.ofcom.org.uk/siteassets/resources/documents/research-and-data/tv-radio-and-on-demand-research/tv-research/news/news-consumption-2024/news-consumption-in-the-uk-2024-report.pdf?v=379621>

<sup>5</sup> E Brogi, G De Gregorio "From the code of practice to the code of conduct? Navigating the future challenges of disinformation regulation" (2024) 16 (1) *Journal of Media Law* 38-46, 40-41; see also the EU Report of the independent High level Group on fake news and online disinformation (2018) [https://ec.europa.eu/newsroom/dae/document.cfm?doc\\_id=50271](https://ec.europa.eu/newsroom/dae/document.cfm?doc_id=50271)

“Twitter files” of 2022<sup>6</sup>, allegations by the right-wing movement that Twitter, before the *Musk* take-over, had been biased in their content moderation towards left-wing views. International journalists such as the Philippine Nobel Peace Prize winner Maria Ressa who has risked her life in uncovering government wrongdoing in her country and who is a staunch supporter of free speech has stated in response to *Zuckerberg’s* announcement that “extremely dangerous times” lied ahead for journalism, democracy and social media users and that a “world without facts” was a “world that’s right for a dictator”.<sup>7</sup> These comments emphasize the need for (social) media organisation’s independence from government and the need for news being based on facts.

This Article examines the connections between disinformation, media regulation, the power of social media, sovereignty and jurisdiction. It proceeds by first presenting two brief case-studies on how disinformation can destabilize democracy, followed by a brief, comparative analysis of the regulation of disinformation in the EU and the US and linking this to questions of jurisdiction and sovereignty. It argues that the concept of sovereignty should ground a principle under international law that states should be able to fight disinformation on big tech platforms as an assertion of their right to govern.

## 2. Disinformation case-studies

### 2.1 What is disinformation?

Academic literature usually makes a distinction between misinformation and disinformation. Disinformation is the intentional spreading of factually incorrect, false news with the intention to mislead the recipients. By contrast, misinformation is false news which was not intentionally spread as such, and the falsity may be due to a honest, but mistaken belief.<sup>8</sup> For the purposes of this article the distinction is not important and I use the term disinformation to mean information presented as news based on false facts and misleading in its message.

### 2.2 Two case-studies

The Article uses two illustrative disinformation case-studies in the context of the summer 2024 riots in the UK<sup>9</sup> and the US 2024 elections.

In July 2024 three little girls were brutally murdered and several others seriously injured by a loner in a mass stabbing incident during a dance class. Immediately after the attack a number of social media posts falsely claimed that the attacker was a Muslim asylum seeker and illegal migrant.<sup>10</sup> This disinformation sparked a series of violent, far-right wing, anti-immigration riots across the UK targeting mosques, shops and immigration hostels with

---

<sup>6</sup> TRT- The Newsmakers (14. December 2022) <https://www.youtube.com/watch?v=sEpIPLzRkZE>

<sup>7</sup> The Guardian (8. January 2025) <https://www.theguardian.com/world/2025/jan/08/facebook-end-factchecking-nobel-peace-prize-winner-maria-ressa>, Rappler (9. January 2025) <https://www.rappler.com/technology/maria-ressa-profit-over-safety-meta-ends-fact-check-program-united-states/>

<sup>8</sup> C Tan The Curious Case of Regulating False News on Google (2022) 46 *Computer Law and Security Review* 1-14, 2; European Commission, Code of Practice on Disinformation (European Commission, September 2022); N Bontridder, Y Pouillet “The role of artificial intelligence in disinformation” (2021) *Data & Policy* doi:10.1017/dap.2021.20 e32-2

<sup>9</sup> The Guardian (22. October 2024) <https://www.theguardian.com/media/2024/oct/22/social-media-algorithms-must-be-adjusted-to-prevent-misinformation-ofcom>

<sup>10</sup> BBC Verify (25. October 2024) <https://www.bbc.co.uk/news/articles/c99v90813j5o>

extensive violence, arson and injuries to police officers for several days. The riots led Elon Musk to further stoke the unrests by his claim that “civil war in the UK is now inevitable”.<sup>11</sup> The media regulator in the UK, Ofcom later published a letter stating that there was a clear link between social media disinformation claiming that the stabbings were committed by an asylum seeker and the violent riots which followed.<sup>12</sup>

Disinformation played a role in the outcome of the 2024 US Election<sup>13</sup>. While disinformation may or may not be overwhelming in terms of quantity of information out there on social media it has influenced political opinions and voting behaviour.<sup>14</sup>

Generative AI and deepfakes have increased the risks in respect of disinformation during elections in every country.<sup>15</sup> Generative AI allows for the speedy and easy generation of inauthentic information, as it has the tendency to include false information and/or generate hallucinated facts.<sup>16</sup> Authentic looking deepfake videos are persuasive in the sense of the old adage that “an image speaks a thousand words” and the same applies to AI’s ability to realistically imitate the human voice. A deepfake has been legally defined in the EU’s AI Act as “AI-generated or manipulated image, audio or video content that resembles existing persons, objects, places, entities or events and would falsely appear to a person to be authentic or truthful”.<sup>17</sup>

The use of AI to imitate *President Biden’s* voice in a misleading robocall made to many voters in the New Hampshire primaries exhorting them not to vote is an example of a mysterious deepfake campaign designed to influence voters through false information.<sup>18</sup>

Famous unsubstantiated stories circulated during the US 2024 election were for example: one (repeated by *Trump*) about Haitian immigrants stealing and eating pets<sup>19</sup>, another the misspending by the Federal Emergency Management Agency of hurricane disaster relief funds on “taking in illegal immigrants”<sup>20</sup>, or another a photo that *Kamala Harris* once embraced convicted sex offender Jeffrey Epstein in beach wear, which was posted on “X”

---

<sup>11</sup> The Telegraph (4. August 2024) <https://www.telegraph.co.uk/news/2024/08/04/southport-latest-news-rioting-disorder-arrests-liverpool/>

<sup>12</sup> <https://www.ofcom.org.uk/siteassets/resources/documents/about-ofcom/public-correspondence/2024/letter-from-dame-melanie-dawes-to-the-secretary-of-state-22-october-2024.pdf?v=383693>

<sup>13</sup> I’m not claiming that disinformation was decisive for the outcome or that disinformation was only perpetrated by one group- all I’m saying here is that it happened.

<sup>14</sup> Elaine C Kamarck, Darrell M West *Lies That Kill- A Citizen’s Guide to Disinformation* (Brookings Institution Press 2024) 19-20

<sup>15</sup> N Bontridder, Y Poullet “The role of artificial intelligence in disinformation” (2021) Data & Policy doi:10.1017/dap.2021.20 e32-3; Brookings “How Disinformation Defined the 2024 Election Narrative” (7. November 2024) <https://www.brookings.edu/articles/how-disinformation-defined-the-2024-election-narrative/>

<sup>16</sup> Seb Butcher “Disinformation, Generative AI and Why Our Laws Need Urgent Reform” (2024) 30 (4) *Computer and Telecommunications Law Review* 85-85, 85

<sup>17</sup> Regulation (EU) 2024/1689 of 13. June 2024, Art 3 (60)

<sup>18</sup> BBC News (22. January 2023) <https://www.bbc.co.uk/news/world-us-canada-68064247>

<sup>19</sup> The Guardian (14. September 2024) <https://www.theguardian.com/us-news/2024/sep/14/racist-history-trump-pet-eating-immigrant>

<sup>20</sup> NBC News (4. October 2024) <https://www.nbcnews.com/politics/donald-trump/false-claims-fema-disaster-funds-migrants-pushed-trump-rcna173955>

shortly after she had announced that she would campaign for the Presidency<sup>21</sup>. She has also been falsely and repeatedly maligned as a “prostitute” or as incompetent for the role.<sup>22</sup>

Examples of Russian false information propaganda have also been uncovered<sup>23</sup>, such as the entirely baseless allegations that the Minnesota Governor and vice-presidential candidate in the 2024 elections *Tim Walz* sexually assaulted one of his students, which have been tracked to a Russian disinformation campaign utilising deepfake whistleblower videos.<sup>24</sup> These stories went viral on social media and some of them were repeatedly mentioned or posted by politicians increasing their reach and impact.

It is incidents like the examples mentioned which risk that we are increasingly living in a post-truth society, where truth is replaced by the power of big tech, weakening elections and other democratic processes, and thereby the security of the democratic state, which is precisely what tech disrupters are trying to achieve.

### **2.3 Targeting of content, manipulation, disinformation**

The starting point of this article are the power struggles between social media platforms manipulating citizens globally, by using complex algorithms to target content to citizens, and amplifying that content with a virality never seen before.<sup>25</sup> This mechanism-perfected to maintain engagement and maximise advertising revenue<sup>26</sup> confers enormous influence and thus power to social media companies. There are two fundamental problems with the virality of social media: first that disinformation frequently has greater reach and velocity<sup>27</sup> and secondly that content targeted for maximum engagement frequently appeals to the subliminal part of the human psyche and emotions, which encourages irrationality and thereby, ultimately popular politics. Algorithm based targeting of content based on behavioural profiling manipulates users.

The targeting of users based on psychological profiles and correlations as to their behaviour has come under much criticism in respect of its negative impacts on society, including privacy infringements, surveillance and lack of autonomy. It is for this reason that the EU Digital Services Act mandates that very large social media platforms and very large search engines have to give users the option to use their recommender systems without profiling,

---

<sup>21</sup> Snopes.com (fact-checking site) <https://www.snopes.com/fact-check/kamala-harris-jeffrey-epstein-beach-photo/>

<sup>22</sup> DW Factcheck (11. July 2024) <https://www.dw.com/en/fact-check-what-role-did-disinformation-play-in-the-us-election/a-70729575>

<sup>23</sup> Similar disinformation campaigns are now targeted at German politicians in the February 2025 elections, EDMO (27. January 2025) “Influence Operation Exposed: How Russia Meddles in Germany’s election campaign” <https://edmo.eu/publications/influence-operation-exposed-how-russia-meddles-in-germanys-election-campaign/>

<sup>24</sup> Wired (21. October 2024)

<sup>25</sup> N Bontridder, Y Poullet “The role of artificial intelligence in disinformation” (2021) *Data & Policy* doi:10.1017/dap.2021.20 e32-3

<sup>26</sup> Shoshana Zuboff *Surveillance Capitalism* (Profile Books 2019)

<sup>27</sup> There are a number of studies which show that disinformation travels faster and leads to more engagement, see for example S Vosoughi, D Roy, S Aral “The Spread of True and False News Online” (2020) 359 *Science* 1146-1151; American Psychology Association Consensus Statement “Using Psychological Science to Understand and Fight Health Misinformation” Report 23. November 2023, explaining why we spread disinformation and funded by the Centers for Disease Control and Prevention of the U.S. Department of Health and Human Services <https://www.apa.org/pubs/reports/misinformation-consensus-statement.pdf>

for example a paid subscription option.<sup>28</sup> Moreover it obliges all service providers to be transparent about their recommender systems and the parameters used.<sup>29</sup> In the US too there is evidence that the law on the divestiture of TikTok was motivated by concerns how a foreign adversary (China) would manipulate the 170 million Tik Tok users in the US through recommender systems targeting content to them.<sup>30</sup>

Disinformation amplified by powerful algorithms has divided societies and challenges rational decision-making and free elections.<sup>31</sup> This state of affairs implies an urgent need to rethink media regulation in order to safeguard the sovereignty of democratic states.

### 3. Content Regulation and the Law on Disinformation

This paper compares the interface between platform regulation and national laws (comparing the EU and US) on disinformation.

#### 3.1 US Law, Content Regulation, Disinformation

In the US tradition, the Constitution largely prohibits the regulation of content (with exceptions) by the state.

Content decisions have been left to private media, including online platforms. Private media are in charge of setting their own content policies and regulating speech through content moderation according to standards described in terms and conditions. Online services are not liable for restricting access to “obscene, lewd, lascivious, filthy, excessively violent, harassing, or otherwise objectionable information” provided by another person, if the provider decides to moderate content on their own initiative (as opposed to being compelled by the law).<sup>32</sup>

Freedom of speech is guaranteed through competition of many different media outlets between which users can *choose according to their own political convictions*. The essence of free speech has been described in *Tik Tok v Garland* as “the principle that each person should decide for himself or herself the ideas and beliefs deserving of expression, consideration, and adherence.”<sup>33</sup>

Harmful speech is reined in through counter-speech, and the notion of the free marketplace of ideas. The theory here is that many viewpoints are juxtaposed and debated, so that ultimately reason and common sense will prevail.<sup>34</sup> Therefore, speech restriction by government is only allowed, if it is absolutely necessary to prevent a serious harm which

---

<sup>28</sup> Art 38, Regulation (EU) 2022/2065 of the European Parliament and of the Council of 19 October 2022 on a Single Market For Digital Services

<sup>29</sup> Arts 27, 40 (3): “explain the design, the logic, the functioning and the testing of their algorithmic systems, including their recommender systems” to the Commission/Digital Services Co-ordinators of the Member States (very large online services only).

<sup>30</sup> *TikTok Inc. v. Garland* Supreme Court of the United States 17. January 2025, p.9

<sup>31</sup> Tim Clement-Jones *Living With the Algorithm- Servant or Master?* (Unicorn Publishing Group 2024) 35-46

<sup>32</sup> 47 US Code §230 (c) (2) (A)

<sup>33</sup> Quoting *Turner Broadcasting System, Inc. v. FCC*, 512 U.S. 622, 641, 114 S.Ct. 2445, 129 L.Ed.2d 497 (1994) on p.5

<sup>34</sup> But it is precisely this mechanism which the behavioural profiling and the individualised micro-targeting of content undermines through echo chambers, and why the “marketplace of ideas” no longer functions see [x-ref](#)

cannot be prevented in any other way ("strict scrutiny")<sup>35</sup>. Under the strict scrutiny test, the government must show that the law is narrow tailored, that it advances a compelling governmental interest and that it is the least restrictive means of implementing this interest, *i.e.* it is necessary in a strict sense.<sup>36</sup> This necessity test is almost impossible to prove.<sup>37</sup> Thus, strict scrutiny means that most harmful speech (such as hate speech) cannot be criminalised by federal and state criminal law. But even criminal speech which is prohibited in conformity with the 1<sup>st</sup> Amendment (such as explicit and realistic child pornography<sup>38</sup>), does not trigger any liability for social media providers hosting this speech. The person who is responsible for the speech (the speaker, the editor) may be criminally liable, but interactive computer services have *absolute immunity from civil liability and state criminal laws* under Section 230 Communications Decency Act 1996<sup>39</sup> and need not take any measures (such as blocking or filtering of content) to prevent the illegal speech from spreading.<sup>40</sup>

This raises the question whether the algorithms actively driving content to users based on their behavioural profiling defeat section 230 CDA. An argument could be made here that social media platforms are not merely passively hosting user provided content, but actively target users with illegal content or particular items of disinformation based on the users' profile and the algorithms calculating that this specific user is *likely to engage with this particular item of content*. This raises the question whether the service goes beyond just being an interactive computer service.<sup>41</sup> Interactive computer services are "any information service, system, or access software provider that provides or enables computer access by multiple users to a computer server, including specifically a service or system that provides access to the Internet and such systems operated or services offered by libraries or educational institutions".<sup>42</sup> This includes all information and computer access systems, as well as hosting services. But arguably the algorithms targeting content specifically based on user profiling are new and may not have been envisaged in 1996- they are not mere information services and access systems.<sup>43</sup> This issue was relevant in two 2023 US Supreme Court cases, but ultimately has not yet been decided by the US Supreme Court. It was effectively side-stepped in *Twitter v Taamneh*<sup>44</sup> and *Gonzalez v Google*<sup>45</sup> where claims that Twitter or Youtube, respectively, were aiding and abetting terrorist acts through enabling terrorist radicalisation were unsuccessful. The decision was reached on the basis that social media platforms' activities did not amount to aiding and abetting so that the Court did not even reach the question of immunity from criminal liability under section 230 CDA. The minimalist ruling has been criticized on the basis that it is a postponement of this important

---

<sup>35</sup> [ACLU v Reno caselaw](#)

<sup>36</sup> *Reed v. Town of Gilbert*, 576 U.S. 155, 163, 135 S.Ct. 2218, 192 L.Ed.2d 236 (2015)

<sup>37</sup> [Ref cases](#)

<sup>38</sup> [refs](#)

<sup>39</sup> 47 US Code §230 (e) (3)

<sup>40</sup> This immunity, however, does not apply to Federal criminal law, 47 US Code §230 (e) (1). If a particular post implicates Federal criminal law, online service providers could theoretically be liable, if they do not take action on being notified.

<sup>41</sup> Casenote Review "Communications Decency Act of 1996" (2023) 137 *Harvard Law Review* 400-409, 405

<sup>42</sup> 47 US Code §230 (5) (f)

<sup>43</sup> [Ref lit and cases](#)

<sup>44</sup> 598 U.S. 471

<sup>45</sup> 598 U.S. 617

question.<sup>46</sup> So while it is clear that online providers are not liable for illegal content, and that content moderation is largely left to the providers themselves, is state legislation which interferes with service providers' content moderation compliant with the 1<sup>st</sup> Amendment? There are two variants to this question, namely first whether (state) legislation can oblige online service to carry out certain types of content moderation and secondly whether (state) legislation can prohibit service providers from moderating certain types of content on their own private initiative.

As to the first variant, there are a number of current legal challenges concerning state laws *mandating content moderation* in respect of children's online safety<sup>47</sup>, i.e. online safety legislation vaguely similar to that found in the EU and UK, if not that restrictive.<sup>48</sup> As to the second variant, legal challenges concerning laws *restricting private platform's content moderation*<sup>49</sup> have reached the US Supreme Court. Broadly generalising the law<sup>50</sup>, currently the position seems to be that both laws mandating content moderation and laws restricting content moderation are unconstitutional, unless they meet the strict scrutiny test<sup>51</sup>.

### 3.1.1 Online Safety Laws Mandating Content Moderation

Several states are in the process of enacting online child safety legislation, termed "Kids Code Bills" or "Age-Appropriate Design Codes".<sup>52</sup> California led the way in 2021 with its Age-Appropriate Design Code Act. These Bills have been successfully opposed by big tech companies though lobbying and First Amendment challenges in the courts. Net Choice<sup>53</sup> is bringing constitutional challenges in the courts, which are resulting in preliminary injunctions against the enforcement of these Acts in several states.<sup>54</sup> This litigation is ongoing at the time of writing. While the take-down of disinformation is not related to child safety, the First Amendment challenges, if successful, make it extremely unlikely that states will successfully enact more restrictive online safety legislation, dealing with child safety, let alone disinformation.

### 3.1.2 Laws Limiting Content Moderation

*Moody v NetChoice* arose out of out of a partisan Republican allegation that social media companies (Google, Meta, Twitter) had been moderating content in an imbalanced way, biased against views on the right and their deletion of accounts of politicians.<sup>55</sup> In response,

---

<sup>46</sup> Casenote Review "Communications Decency Act of 1996" (2023) 137 Harvard Law Review 400-409, 409

<sup>47</sup> *NetChoice v Bonta* 113 F.4 1101 (9<sup>th</sup> Cir. California August 2024), *NetChoice v Reyes* 2024 WL 4135626 (US District Court D.Utah September 2024; appeal to 10. Cir pending)

<sup>48</sup> Discussed below, ref

<sup>49</sup> *Moody v NetChoice; Paxton v Net Choice* 144 S.Ct. 2383 (US Supreme Court, 1. July 2024)

<sup>50</sup> This broad characterization is inaccurate, since the cases are currently making their way through the courts and the assessment depends on the precise provisions in state legislation challenged, but the generalisation helps to conceptualise the broader direction of travel. More detail is provided below.

<sup>51</sup> See FN ref 37

<sup>52</sup> California, Maryland, Vermont, Minnesota, Hawaii, Illinois, New Mexico, Nevada, South Carolina, Utah

<sup>53</sup> Net Choice is the trade association for the big tech companies, including Alphabet/Google, Meta, X, Snap and Pinterest, tasked with minimising regulation of their activities, see <https://netchoice.org/about/>

<sup>54</sup> *NetChoice v Bonta* 113 F.4 1101 (9<sup>th</sup> Cir. California August 2024), *NetChoice v Reyes* 2024 WL 4135626 (US District Court D.Utah September 2024; appeal to 10. Cir pending)

<sup>55</sup> BBC News (date) <https://www.bbc.co.uk/news/world-us-canada-55597840>

Florida and Texas passed state legislation prohibiting the deplatforming of politicians and viewpoint discrimination, thereby restricting platforms' mechanisms of content moderation. NetChoice brought a legal challenge against both laws in the respective US District Courts<sup>56</sup> both of which issued a preliminary injunction on the basis that these laws triggered strict scrutiny under the 1<sup>st</sup> Amendment which neither law satisfied. These conflicting decisions were eventually joined before the US Supreme Court in *Moody v NetChoice*<sup>57</sup>, which unanimously vacated the appeal judgments and remanded the cases back to the lower courts<sup>58</sup>, without deciding the issues fully. The US Supreme Court held that social media platforms produced their own compilations of expressions which were protected under the 1<sup>st</sup> Amendment, so that the prohibition of viewpoint discrimination might engage the *free speech rights of the social media companies* (as opposed to the users of the platforms).<sup>59</sup> Thus, even though in practice the platforms removed very little speech and were not exercising editorial functions as such, the fact that they compiled the expressions of users on their platform through algorithms gave the social media companies 1<sup>st</sup> Amendment protection.<sup>60</sup> The Supreme Court held that it was not the government's role to balance the free marketplace of ideas.<sup>61</sup> But ultimately it is for the lower courts to examine whether these Laws are in breach of the 1<sup>st</sup> Amendment<sup>62</sup> and commentators have pointed out that final resolution of these issues might be years away.<sup>63</sup>

These two aspects of the US Supreme Court view, namely that algorithmic content curation through micro-targeting of users and that governments does not interfere with the marketplace of ideas in the name of pluralism and diversification of views is precisely what distinguishes US law from EU law, where it is part of the state's role to ensure pluralism of ideas and where algorithmic targeting is seen as a threat to freedom of expression, and in particular the right to receiving information.<sup>64</sup>

### 3.1.3 Informal Nudging by the State to Encourage Online Services to Limit the Spread of Disinformation

Another tension between freedom of speech and content moderation was raised before the Supreme Court in 2024 in *Murthy v Missouri*<sup>65</sup>. The allegation was that the Biden administration had directed the social media companies to take down disinformation. In the wake of these allegations, the then Attorney Generals of Missouri and Louisiana and right-wing user groups applied for an injunction mandating the Biden administration not to urge the social media companies to act on disinformation. The injunction was initially granted<sup>66</sup>,

---

<sup>56</sup> *Moody v NetChoice* US District Court for the Northern District of Florida 546 F.Supp.3d 1082 and *NetChoice v Paxton* US District Court for the Western District of Texas 573 F.Supp.3d 1092

<sup>57</sup> *Moody v NetChoice* 144 S.Ct. 2383; 219 L.Ed.2d 1075 (1. July 2024)

<sup>58</sup> Justice Kagan gave the Opinion of the Court, with four concurring opinions **double-check**

<sup>59</sup> At 2393, this was criticised by the concurring judgment of Justice Alito at 2431-2432

<sup>60</sup> At 2402

<sup>61</sup> At 2402-2403

<sup>62</sup> At 2409

<sup>63</sup> E Goldman, "Speech Nirvanas on the Internet: an Analysis of the U.S. Supreme Court's *Moody v NetChoice* Decision" 2024 *Cato Supreme Court Review* 125-155, 126

<sup>64</sup> **Cross-ref and compare later discussion**

<sup>65</sup> 144 S.Ct. 1972; 219 L.Ed.2d 604

<sup>66</sup> 680 F.Supp.3d 630; US District Court for the Western District of Louisiana



then narrowed on appeal to the Fifth Circuit<sup>67</sup>, and finally stayed with an appeal to the US Supreme Court<sup>68</sup>. The US Supreme Court<sup>69</sup> held that the injunction was invalid for the reason that the claimants had no standing under Article III of the US Constitution<sup>70</sup> which only allowed judicial review where this was “necessary to *redress or prevent actual or imminently threatened injury* to persons caused by official violation of law” (my emphasis).<sup>71</sup> Since the claimants relied on *past* allegations of government interference and could not prove a chain of *causation* between the administrations’ communications and the take-down or blocking of material there was no *actual or imminent threat* to their freedom of speech justifying the imposition of an injunction.<sup>72</sup> The Court emphasized that content moderation by social media companies was not as such unconstitutional and it found that they moderate content, so it could not be proven that content was removed *because of* state interference.<sup>73</sup>

The Supreme Court held in *Murthy* that a degree of political pressure can be applied by the state, as long as the ultimate responsibility for content moderation is free from actual government interference. The *Murthy* decision is significant in that it disallowed an injunction against Government influencing social media content moderation practices to a degree, including the *flagging of posts constituting disinformation information*.

President Trump in his first days in office has passed an Executive Order “Restoring Freedom of Speech and Ending Federal Censorship”<sup>74</sup>. Already in 2022, he vowed to “shatter the left-wing censorship regime”.<sup>75</sup> This Executive Order, in its first part, criticised the Biden administration for trampling on

“free speech rights by censoring Americans’ speech on online platforms, often by exerting substantial coercive pressure on third parties, such as social media companies, to moderate, deplatform, or otherwise suppress speech that the Federal Government did not approve”.<sup>76</sup>

While this first part gives context, it does not have retroactive effect, it simply criticises what has happened in the past. In the second part the Executive Order directs the federal (Trump) administration to “ensure that no Federal Government officer, employee, or agent engages in or facilitates any conduct that would unconstitutionally abridge the free speech of any American citizen”<sup>77</sup>. While the political intention of this exhortation is clear, its legal significance is at least doubtful. Either a federal act is an infringement of free speech, then this would be illegal under the 1<sup>st</sup> Amendment anyway, or it does not engage free speech

---

<sup>67</sup> 83 F.4th 350

<sup>68</sup> 144 S.Ct. 7

<sup>69</sup> Opinion written by Justice Amy Coney Barrett

<sup>70</sup> U.S. Const. art. 3, § 2, cl. 1; Fn 65 at 1985-97

<sup>71</sup> Citing *Clapper v. Amnesty Int’l USA*, 568 U.S. 398, 409, 133 S.Ct. 1138, 185 L.Ed.2d 264

<sup>72</sup> At 1987

<sup>73</sup> Fn 65 at 1992-96

<sup>74</sup> January 20., 2025 see <https://www.whitehouse.gov/presidential-actions/2025/01/restoring-freedom-of-speech-and-ending-federal-censorship/>

<sup>75</sup> Sara Dorn “Trump Vows To Dismantle ‘Censorship Cartel’ If He’s Re-Elected—An Apparent Nod To Musk’s ‘Twitter Files’ Release” (Forbes 15. December 2022) <https://www.forbes.com/sites/saradorn/2022/12/15/trump-vows-to-dismantle-censorship-cartel-if-hes-re-elected-an-apparent-nod-to-musks-twitter-files-release/>

<sup>76</sup> Clause 1 Purpose

<sup>77</sup> Clause 2 (b) and 3 (a)

rights, in which case the clause is not applicable. It is difficult to see how this Executive Order adds to 1<sup>st</sup> Amendment Rights or their interpretation by the courts. More significant are Clauses 2 (c) and 3 (a) which state that no federal resources must be deployed for fighting disinformation online which means that no federal agency is able to use resources and manpower to liaise with social media companies and search engines to carry out that task.<sup>78</sup> Even more concerningly, these Clauses may mean that any third party, independent of the federal government, but supported by federal funding, such as non-profits, universities and colleges may have their funding withdrawn if they engage in activities such as content flagging or fact-checking, or even research which may impact the free flow of disinformation.<sup>79</sup>

Furthermore Clauses 2 (d) and 3 (a) provide that it is the policy of the Trump administration to “identify and take appropriate action to correct past misconduct by the Federal Government related to censorship of protected speech”. This raises the question whether officers who were fighting disinformation (for example during the Pandemic) will be dismissed or subjected to some disciplinary action in their employment. Such disciplinary action would be retroactive and raises questions about its legality. The wording is “to correct past misconduct”- this is an impossible task as disinformation repressed in the past cannot be effectively corrected, possibly years later. The correction may be a right of reply of the person whose speech had been suppressed or some duty to correct the record- but how effective this would be is the question- a real Don Quixote remedy. Therefore, ultimately it is not clear what the correction entails. However, for the purposes of remedial action, the Executive Order mandates that the Attorney Generals should investigate and prepare a Report submitted to the President with recommendations for remedial action.<sup>80</sup> The Executive Order makes clear that it does not confer any rights of action to persons who claim that their speech has been curtailed illegally.<sup>81</sup>

### **3.2 European Laws: Media Regulation in the Service of Freedom of Expression**

On the other side of the Atlantic, in democratic European countries, speech regulation is an accepted form of media regulation and involves the careful balancing of freedom of expression with harm stemming from speech. This is reflected in the European Convention of Human Rights, which in Article 10 protects the right to freedom of expression which is conceived as both a right to speak as well as a right to receive information<sup>82</sup>, and arguably it is the right to receive truthful information which is engaged by the spreading of disinformation, or the state not preventing the spread of disinformation. But the right to freedom of expression is constrained where this is absolutely necessary to protect the rights

---

<sup>78</sup> “ensure that no taxpayer resources are used to engage in or facilitate any conduct that would unconstitutionally abridge the free speech of any American citizen.”

<sup>79</sup> Sara Dorn “Trump Vows To Dismantle ‘Censorship Cartel’ If He’s Re-Elected—An Apparent Nod To Musk’s ‘Twitter Files’ Release” (Forbes 15. December 2022)  
<https://www.forbes.com/sites/saradorn/2022/12/15/trump-vows-to-dismantle-censorship-cartel-if-hes-re-elected-an-apparent-nod-to-musks-twitter-files-release/>

<sup>80</sup> Clause 3 (b)

<sup>81</sup> Clause 4 (c)

<sup>82</sup> Art 10 (1) “Everyone has the right to freedom of expression. This right shall include freedom to hold opinions and to receive and impart information and ideas without interference by public authority and regardless of frontiers. This Article shall not prevent States from requiring the licensing of broadcasting, television or cinema enterprises.”

of others as prescribed by the law, including laws covering false and harmful speech.<sup>83</sup> European media law, by way of generalisation is characterised by arms-length regulation through independent<sup>84</sup> regulators<sup>85</sup> and ultimately independent courts who carry out this balance. Media regulation in Europe is broadly characterised by the goal to ensure plurality of media voices.<sup>86</sup> Additionally it has been argued that users of media services should be exposed to a number of different voices and viewpoints to ensure media plurality is effective and that this should be one of the goals of regulation.<sup>87</sup> *Koltay* has pointed out that mass media should be a forum enabling public debate, but in practice commercial interests have negatively impacted this function, such that the media “treat the public less as a participant in public debate and more as a revenue-generating audience”.<sup>88</sup>

In view of this conception of media law, it is only logical that European laws have instituted regulation of social media services in the shape of the Digital Services Act (EU)<sup>89</sup> and the Online Services Act (UK<sup>90</sup>), which both regulate social media platforms and search engines.

In the EU media regulation is not focused on restricting illegal or harmful content, but rather on the dynamics of disinformation, including the algorithms which amplify and target content and the responsibility of gatekeepers.<sup>91</sup>

### 3.2.1 The EU Digital Services Act

The Digital Services Act<sup>92</sup> imposes obligations on social media companies to moderate content based on the notification of illegal content (including by trusted flaggers such as the police/state authorities or victim organizations whose notifications should be prioritised<sup>93</sup>). Furthermore, online service providers have to restrict illegal content and notify illegal content, which they are aware of, to the (police) authorities of the relevant Member State. Online service providers have to operate systems which allow users to notify illegal content, and instituting action to prevent continued access (“notice and take-down”, “notice and

---

<sup>83</sup> Art 10 (2) “The exercise of these freedoms, since it carries with it duties and responsibilities, may be subject to such formalities, conditions, restrictions or penalties as are prescribed by law and are necessary in a democratic society, in the interests of national security, territorial integrity or public safety, for the prevention of disorder or crime, for the protection of health or morals, for the protection of the reputation or rights of others, for preventing the disclosure of information received in confidence, or for maintaining the authority and impartiality of the judiciary.”

<sup>84</sup> Independent of the government, albeit that practice may not always follow the principle (which is another discussion for another article)

<sup>85</sup> This is also reflected in the Digital Services Act, Art 52 (1) and (2): “When carrying out their tasks and exercising their powers in accordance with this Regulation, the Digital Services Coordinators shall act with complete independence. They shall remain free from any external influence, whether direct or indirect, and shall neither seek nor take instructions from any other public authority or any private party.”

<sup>86</sup> A Koltay *Media Freedom and the Law* (Routledge 2025) 141, 158, 161

<sup>87</sup> P M Napoli “Exposure Diversity Reconsidered” (2011) 1 *Journal of Information Policy* 246-259

<sup>88</sup> FN 86, 162

<sup>89</sup> In force since 17. February 2024

<sup>90</sup> The United Kingdom left the European Union after a referendum in 2016 on 31. January 2020. The referendum campaign to leave was influenced by serious misinformation, such as the famous campaign bus, which claimed that the state would save £350 million which could go to the National Health Service, BBC News (16. January 2018) <https://www.bbc.co.uk/news/uk-42698981>

<sup>91</sup> E Brogi, G De Gregorio “From the code of practice to the code of conduct? Navigating the future challenges of disinformation regulation“ (2024) 16 (1) *Journal of Media Law* 38-46, 39

<sup>92</sup> Regulation (EU) 2022/2065 of the European Parliament and of the Council of 19 October 2022 on a Single Market For Digital Services; came into force on 17. February 2024, Art 93

<sup>93</sup> DSA Art 22

action”).<sup>94</sup> Service providers<sup>95</sup> have to provide complaints systems, allowing both, for complaints that content has been unfairly restricted or, the reverse, where a notification of illegal content has been ignored.<sup>96</sup> Complainants must be provided with independent Alternative Dispute Resolution, if their complaint has not been solved to their satisfaction.<sup>97</sup> Complainants also have a right to complain to their national regulatory authority.<sup>98</sup> There are also obligations to prepare Transparency Reports about content moderation.<sup>99</sup> Very large online services<sup>100</sup> (both social media platforms and search engines) have to comply with additional obligations<sup>101</sup>. These companies have to carry out systemic risk assessments (including risks stemming from algorithms) as to harms arising from illegal content or illegal activities on their services<sup>102</sup> and mitigate these risks proactively.<sup>103</sup> The risk assessments must include the risks for “inauthentic use”, ie the spreading of disinformation.<sup>104</sup> The risk mitigation requirements are extensive and include a whole range of measures in Art. 35 (1), including adapting the design and functioning of the service, adapting and increasing content moderation, adapting and testing algorithms, including content recommender systems, adapting advertising systems, demonetizing content, awareness raising measures such as explaining why content is being targeted or fact-checks in respect of false information, measures protecting children such as age-verification, parental controls or support measures (for example for mental health problems or eating disorders).<sup>105</sup> Deepfakes must be prominently and clearly marked as false or inauthentic (for example in respect of humour or parody) to ensure that they are not misleading. There must be an interface<sup>106</sup> for users enabling them to mark up inauthentic content.<sup>107</sup> The Digital Services Act additionally requires specific measure to respond to crises in public security or public health<sup>108</sup>, to deal with the negative consequences of situations such as rioting, a major

---

<sup>94</sup> Take-down is only one option- other action can include restricting accounts, deprioritising or demonetizing content etc, Arts 10, 16-17 DSA

<sup>95</sup> Other than micro- and small enterprises, Art 19

<sup>96</sup> Art 20 DSA

<sup>97</sup> Art 21 DSA; I have suggested the introduction of ADR for such disputes already 23 years ago: J. Hornle “Internet Service Provider Liability – Let’s Not Play Piggy in the Middle” (2002) 7(3) *Communications Law* 85–89

<sup>98</sup> Art 53

<sup>99</sup> Art 24 and Art 42 (for very large online services this includes details about human content moderators, broken down by language and accuracy) DSA

<sup>100</sup> Defined as having more than 45 million users, Art 33 (1) DSA

<sup>101</sup> The European Commission has designated the very large online platform services and very large online search providers, see <https://digital-strategy.ec.europa.eu/en/library/designation-decisions-first-set-very-large-online-platforms-vlops-and-very-large-online-search> (20. December 2023)

<sup>102</sup> The risks are not limited to illegality, and include risks to fundamental rights such as users’ freedom of expression, viewpoint plurality, civic discourse and electoral processes, and public security and negative effects on gender based violence; negative effects on children and public health, Art 34 (1)

<sup>103</sup> Arts 34-35 DSA

<sup>104</sup> Art 34 (1) DSA

<sup>105</sup> Art 35 (1) DSA

<sup>106</sup> This can take the **shape of Community Notes**, see further the discussion **x-ref**

<sup>107</sup> Art 35 (1) (k) DSA: “ensuring that an item of information, whether it constitutes a generated or manipulated image, audio or video that appreciably resembles existing persons, objects, places or other entities or events and falsely appears to a person to be authentic or truthful is distinguishable through prominent markings when presented on their online interfaces, and, in addition, providing an easy to use functionality which enables recipients of the service to indicate such information. “

<sup>108</sup> A “crisis shall be deemed to have occurred where extraordinary circumstances lead to a serious threat to public security or public health in the Union or in significant parts of it”, Art 36 (2)

terrorist incident or a pandemic further fanned, for example, by social media disinformation.<sup>109</sup> The Commission makes the decision whether such a crisis has arisen to kick the crisis management in action.<sup>110</sup> The crisis provisions require the service providers to conduct risk assessments specific to the crisis situation and take measures in dialogue with the European Commission, taking into account civil liberties, and then report on the measures taken to the Commission.<sup>111</sup> In terms of the crisis fighting measures required, the DSA refers back to Article 35 (1) and states that the restrictive measures have to be proportionate to the public security or public health risks<sup>112</sup>. Furthermore, the Commission is drawing up crisis management protocols in advance and in co-operation with very large (and other) online service providers.<sup>113</sup> The measures could involve restricting the accounts of certain groups or slowing down the amplification of content, speeding up notification systems for illegal content and beefing up the concomitant content moderation, or imposing fact-checking measures. The online service providers themselves determine the actual, specific measures to be adopted<sup>114</sup> in a “dialogue” with the European Commission<sup>115</sup>. The DSA envisages that the Commission effectively takes the role of protector of civil liberties (safeguarding privacy and freedom of expression), by checking that the measures are proportionate, subject to safeguards and time-limited.<sup>116</sup> The Commission in turn has to report to and consult the new European Board for Digital Services<sup>117</sup> in order to ensure that it does not fail in this role.<sup>118</sup>

As in the US the DSA leaves the technical decision on the measures and their implementation to online service providers. However, the initiation of the crisis management procedure is left the European Commission who also has to ensure it complies with civil liberties. This conception of the state (and here a supranational organisation) as the guarantor of civil liberties is fundamentally different to the US, where this is an alien notion. The European Commission is the main enforcer of the DSA as regards the obligations of very large online service providers, with co-operation of the Member States’ regulatory authorities.<sup>119</sup>

Moreover, the DSA has stringent<sup>120</sup> and extensive enforcement provisions, including the requirement to have independent, annual compliance audits<sup>121</sup>, a compliance function<sup>122</sup>, investigatory powers, including against third parties, such as the independent auditors<sup>123</sup>,

---

<sup>109</sup> Art 36 DSA

<sup>110</sup> Art 36 (1)

<sup>111</sup> Art 36 (1) (a), (b) and (c)

<sup>112</sup> “specific, effective and proportionate measures” Art 36 (1) (b)

<sup>113</sup> Art 48 (2)

<sup>114</sup> Art 36 (5)

<sup>115</sup> Art 36 (4), (6)

<sup>116</sup> Art 36 (1) (b), (3), (8) (b)

<sup>117</sup> Art 61 composed of the so-called Digital Services Co-ordinators, ie the national regulatory authorities

<sup>118</sup> Art 36 (4) (c), (7), (8), (10)

<sup>119</sup> Arts 56-57; this avoids the “Irish problem”, ie the fact that most very large online service providers are headquartered in Ireland, and the Irish regulators having insufficient resources and willpower to effectively enforce the law [ref](#); Arts 65-66 on the Commission’s enforcement procedure.

<sup>120</sup> Some may say “draconian” enforcement powers, which extend to third parties holding relevant information, Art 51 (2)

<sup>121</sup> Art 37 (very large online services)

<sup>122</sup> Art 41 (very large online services)

<sup>123</sup> Art 51 (1) (national regulators); Arts 67-69 (European Commission)

and significant fines for non-compliance, including daily penalties<sup>124</sup>. Users can claim compensations from service providers for failure to comply with the DSA.<sup>125</sup> As a measure of last resort, the national regulators of the Member States may order the management of online services to set out an “action plan” to avert identified risks, and if this is not forthcoming, the relevant service may be temporarily blocked by a court order.<sup>126</sup> The court proceedings for these wholesale access blocking orders must allow *amicus curae* representations.<sup>127</sup> All enforcement measures, including the blocking order must weigh up the seriousness of the harms to be prevented with the restrictions on freedom of expression and the feasibility and impact on the business of online services<sup>128</sup>. In other words, enforcement measures have to comply with the principle of proportionality.<sup>129</sup>

The DSA additionally mandates access by vetted researchers to online service providers’ data and systems to enable them to independently research the systemic risks posed by very large social media services and search engines.<sup>130</sup> It also encourages co-regulatory Codes of Conduct.<sup>131</sup>

### 3.2.2 Voluntary Code of Practice on Disinformation

One example of a Code, which precedes the DSA is the EU co-regulatory Code of Practice on Disinformation.<sup>132</sup> This Code contains *voluntary* commitments negotiated with online service platforms to fight disinformation through a number of actions. In particular, there are seven main areas on which the Code focuses, namely: identifying and demonetizing disinformation including rules on ad placements, disclosures related to political advertising, agreeing a definition of impermissible manipulative behaviours, enhanced co-operation with fact-checkers, user empowerment and digital literacy, guaranteeing that researchers are having access to data and, the creation of a new transparency centre providing data about disinformation and information sharing. The signatories<sup>133</sup> have entered into commitments in respect of measures with Key Performance Indicators attached in each of these seven areas. The Code led to the establishment of a permanent task force, supported by the European Digital Media Observatory (EDMO). EDMO uses the resources of 14 research hubs covering the entire 27 Member States of the EU. The “hubs” are consortia composed of research institutes and non-profit media and free speech organisations (including fact-checking organisations). Their role is to identify and analyse disinformation campaigns, playing a supporting role in investigative journalism to expose such campaigns, develop tools to analyse the provenance of information, including deepfakes, organising media literacy

---

<sup>124</sup> Art 52 (3)- (4) and Arts 74, 76: a maximum of 6% of annual worldwide turnover in the preceding financial year, or periodic penalties of 5 % of the average daily worldwide turnover

<sup>125</sup> Art 54

<sup>126</sup> Art 51 (3); Arts 75-76, 82 (European Commission in respect of very large online services)

<sup>127</sup> Art 51 (3) Second sentence

<sup>128</sup> Art 51 (5): “the economic, technical and operational capacity of the provider of the intermediary services”

<sup>129</sup> Art 51 (3) Second sentence, and (5); Art 36 (1) (b), (3), (8) (b)

<sup>130</sup> Art 40

<sup>131</sup> Arts 45, 46, 47

<sup>132</sup> “Strengthened Code of Practice on Disinformation” of 16. June 2022, which built on its predecessor, the 2018 self-regulatory Code of Conduct, see <https://digital-strategy.ec.europa.eu/en/library/2022-strengthened-code-practice-disinformation>

<sup>133</sup> The original signatories of the voluntary Code include Google, Meta, Microsoft and Tik-Tok; unsurprisingly “X” is not part of this system (16. June 2022), see <https://digital-strategy.ec.europa.eu/en/library/signatories-2022-strengthened-code-practice-disinformation>

activities to citizens in their region, monitoring the digital information system and feeding this intelligence to national media regulators and other public authorities.<sup>134</sup> The significance of the voluntary Code on Disinformation is that it has created an infrastructure for transparency reporting and framework for research to better understand the structural indicators of disinformation and the impact of measures taken. It has been argued that the metrics should be improved across platforms and that the audit requirements under the DSA, the access for researchers to the data of very large service providers and the Code of Practice on Disinformation should be better connected.<sup>135</sup> The following structural indicators have been suggested to measure the effectiveness of the voluntary Code: prevalence of a disinformation campaign, its sources, audience engagement, disinformation revenue, measuring cooperation and investments in fact-checking, and support for implementing the Code by platforms.<sup>136</sup>

### 3.2.3 EU Regulation of Artificial Intelligence and Disinformation

In addition to the DSA and the Code on Disinformation, the EU Artificial Intelligence Act<sup>137</sup> is relevant to the containment of disinformation. Artificial Intelligence plays a role in dissemination of disinformation and its containment in four respects. First of all, as already pointed out, Generative AI is enabling the creation of disinformation. Secondly, AI is used in the behavioural profiling of social media and search users and the targeting and matching of content. Thirdly the use of deepfakes in disinformation also involves AI. Finally, tools detecting and labelling disinformation are increasingly AI based and their malfunctioning impacts freedom of expression. This raises the question of how the new AI regulation applies to these four uses of AI.

As to Generative AI, this falls into the category of General Purpose AI (GPAI) Models<sup>138</sup>. Where Generative AI models are used to generate information or advertising there is a risk that the content generated contains mistakes and so-called hallucinations<sup>139</sup>, which leads to the spreading of disinformation.<sup>140</sup> Alternatively, AI may be created to purposely spread disinformation, and the Model is trained to generate the most fantastic and outlandish content with a view to increasing engagement and virality of news stories, whether the motive is income generation or ideology. The EU AI Act envisages that the issue of synthetically generated or manipulated information can be addressed by technological means of detection and labelling according to standards laid down in Codes of Conduct.<sup>141</sup>

---

<sup>134</sup> EDMO “About Us” <https://edmo.eu/about-us/edmo-hubs/>

<sup>135</sup> S. Lai, K Yadav “Operational Reporting in Practice: the EU’s Code of Practice on Disinformation” Carnegie Research Article (21. November 2023) <https://carnegieendowment.org/research/2023/11/operational-reporting-in-practice-the-eus-code-of-practice-on-disinformation?lang=en>

<sup>136</sup> E Brogi, G De Gregorio “From the code of practice to the code of conduct? Navigating the future challenges of disinformation regulation“ (2024) 16 (1) Journal of Media Law 38-46, 45

<sup>137</sup> Regulation (EU) 2024/1689 of 13. June 2024

<sup>138</sup> A GPAI Model “displays significant generality and is capable of competently performing a wide range of distinct tasks regardless of the way the model is placed on the market and that can be integrated into a variety of downstream systems or applications”, Art 3 (63) AI Act

<sup>139</sup> This occurs where Large Language Models identify non-existent patterns, see IBM, “What are AI hallucinations?” (1. September 2023) <https://www.ibm.com/think/topics/ai-hallucinations>

<sup>140</sup> The use of Generative AI for the spreading of disinformation is foreseen by the Regulation, see Recital 120

<sup>141</sup> Art 56 and Recitals 135-136

It regulates Generative AI in three ways. First it lays down disclosure requirements on *users*<sup>142</sup> and *providers*<sup>143</sup> of Generative AI.<sup>144</sup> *Users of AI* generating or manipulating text which is published with the purpose of informing the public on matters of public interest (news) must disclose that the text has been artificially generated or manipulated.<sup>145</sup> This disclosure obligation does not apply if there is human review and editorial responsibility.<sup>146</sup> The definition of users does not encompass private, personal users, so that the disclosure obligation also does not apply to private and personal use of Generative AI.<sup>147</sup> *Providers of Generative AI* must ensure that the outputs of the AI system are marked in a machine-readable format and detectable as artificially generated or manipulated. These technical solutions must be effective, interoperable, robust and reliable as far as this is technically feasible.<sup>148</sup>

Secondly it provides specific standards and obligations on Generative AI systems with systemic risks.<sup>149</sup> A systemic risk is defined as “a risk that is specific to the high-impact capabilities of general-purpose AI models, having a significant impact on the [EU] market due to their reach, or due to actual or reasonably foreseeable negative effects on public health, safety, public security, fundamental rights, or the society as a whole, that can be propagated at scale”.<sup>150</sup> Generative AI systems made to disseminate disinformation may fit into this category but the classification of an AI system depends on classification by the EU Commission<sup>151</sup> according to high impact capabilities benchmarks<sup>152</sup>. If an AI model is classified as having systemic risks, such as the wide spreading of disinformation, the provider has to perform “model evaluation in accordance with standardised protocols and tools reflecting the state of the art”, conduct and document adversarial testing, and mitigate the risks identified, such as a risk that the model leads to the spreading of disinformation.<sup>153</sup>

Thirdly it imposes documentation and transparency obligations on all Generative AI systems, including training and testing processes.<sup>154</sup> Providers of General Purpose AI models placed on the market in the EU must appoint a legal representative within the EU<sup>155</sup> to ensure compliance and co-operation with the AI Office.

---

<sup>142</sup> Called deployers of AI, defined in Art 3 (4) AI Act

<sup>143</sup> The legislative term is “providers of an AI system”, Art 3 (3) AI Act: “natural or legal person, public authority, agency or other body that develops an AI system or a general-purpose AI model or that has an AI system or a general-purpose AI model developed and places it on the market or puts the AI system into service under its own name or trademark (...)”

<sup>144</sup> Art 50 (2) AI Act: “Providers of AI systems, including general-purpose AI systems, generating synthetic audio, image, video or text content, shall ensure that the outputs of the AI system are marked in a machine-readable format and detectable as artificially generated or manipulated.”

<sup>145</sup> Art 50 (4) AI Act

<sup>146</sup> Art 50 (4) AI Act

<sup>147</sup> Art 3 (4) AI Act

<sup>148</sup> Art 50 (2) AI Act

<sup>149</sup> Art 55 AI Act

<sup>150</sup> Art 3 (65)

<sup>151</sup> Art 51 (1) (b) AI Act

<sup>152</sup> See Art 51 (1) and Annex XIII AI Act; in vague terms: the reach, size and autonomy of the AI

<sup>153</sup> Art 55 (1) (a)

<sup>154</sup> Art 53 (1) (a) and Annex XI Evaluation

<sup>155</sup> Art 54 (1)



As to deepfakes<sup>156</sup>, the AI Act does not prohibit them, but imposes marking and disclosure requirements, essentially the same as for General Purpose AI.

First, *users*<sup>157</sup> *generating the deepfake content* must disclose that the content/output is artificially generated or manipulated.<sup>158</sup> However, if the content “forms part of an evidently artistic, creative, satirical, fictional or analogous work or programme”, the disclosure obligations are limited in order to ensure that they do not interfere with the display or enjoyment of the work.<sup>159</sup>

Secondly, the *providers of the AI*<sup>160</sup> for generating the inauthentic or false content have to mark the outputs in a machine-readable format, so that its nature becomes detectable by AI.<sup>161</sup> This requirement applies regardless of the nature of the output or its use. Thus, the AI Act mainly relies on technological means to detect deepfakes and these technical methods must be “effective, interoperable, robust and reliable as far as this is technically feasible, taking into account (...) the costs of implementation and the generally acknowledged state of the art, as may be reflected in relevant technical standards.”<sup>162</sup>

### **3.2.4 Immunity from Legal Liability for Online Service Providers in the EU**

As in the US, the EU has legislative provisions giving online services, including social media, immunity for the content they host.<sup>163</sup> But unlike the US, the immunity is based on knowledge, in other words they are immune from civil and criminal liability, unless they have actual or constructive knowledge of the illegal content or activities.<sup>164</sup> EU Member States must not impose obligations on online services to filter *all* content for illegal items.<sup>165</sup> There is clearly a tension between the requirement for very large online service providers to carry out risk assessments and risk mitigation, on the one hand, and the exclusion of a general monitoring obligation, on the other hand.<sup>166</sup>

There are fundamental differences in media regulation between Europe and the USA. In Europe media plurality is achieved by state interference regulating media service providers and balancing commercial freedom with obligations as to plurality and content moderation, for example the prevention of disinformation. In the US media plurality is achieved by encouraging competition and leaving regulation to the marketplace of ideas, including technological innovation. These two approaches reflect how social media services are

---

<sup>156</sup> See FN 17X

<sup>157</sup> See FN 142

<sup>158</sup> Art 50 (4) AI Act

<sup>159</sup> Art 50 (4) AI Act

<sup>160</sup> See FN 143

<sup>161</sup> Art 50 (2) “detectable as artificially generated or manipulated”

<sup>162</sup> Art 50 (2) AI Act

<sup>163</sup> This had been first introduced by E-commerce Directive 2000/31/EC, and is now in the DSA

<sup>164</sup> Art 6 DSA: “the service provider shall not be liable for the information stored (...) on condition that the provider: (a) does not have actual knowledge of illegal activity or illegal content and, as regards claims for damages, is not aware of facts or circumstances from which the illegal activity or illegal content is apparent; or (b) upon obtaining such knowledge or awareness, acts expeditiously to remove or to disable access to the illegal content.”

<sup>165</sup> Art 8 DSA

<sup>166</sup> Ref to cases

regulated and are fundamentally incompatible<sup>167</sup>, hence creating a conflict of law on the internet.

#### 4. How does this relate to jurisdiction and sovereignty?

The differences in the regulation of social media is increasingly leading to conflicts of laws between EU States, the UK and the US. It is here that the concepts of jurisdiction and sovereignty are relevant. The essence of jurisdiction is the authority to state the law, which is literally what this Latin term means. The concept concerns the competence of a state's authorities to pass laws, adjudicate disputes and enforce the laws and it is territorially bounded, or attaches to nationality or domicile of persons.<sup>168</sup> The source of the authority in democratic countries is the people conferring power on their government.<sup>169</sup> Jurisdiction is not limited to a state's territory and is overlapping in the sense that more than one state may assert jurisdiction over the same activity.<sup>170</sup> Furthermore, even if a state asserts jurisdiction, in practice the state may not be able to enforce the law it wishes to apply, so as well as jurisdiction being overlapping, it is frequently incomplete or ineffective.<sup>171</sup> This jurisdictional enforcement challenge impacts a state's sovereignty, as it frequently limits the state's ability to govern activities online.<sup>172</sup>

Sovereignty, then, as a concept of international law describes a state's *power* over its territory, and as a corollary, over its population and the activities located in its territory.<sup>173</sup> It is agnostic as to the source of the power and its legitimacy. It is disputed whether sovereignty is merely a concept to explain statehood, or confers a bundle of rights to states, including the right to non-interference.<sup>174</sup> The concept of state sovereignty has recently come under attack.<sup>175</sup> But while legal systems are still composed of states, statehood is inexorably associated with sovereignty, and to that extent, it is a necessary concept.<sup>176</sup> Arising from the concept of sovereignty is the notion of equality of states under international law regardless of their economic or political power, and that states have a right of non-interference with their domestic affairs, both of which are fundamental tenets of the modern international legal system, even if not always obeyed in practice.<sup>177</sup>

As a matter of *international law* what do these two concepts of jurisdiction as *authority* and sovereignty as *power* tell us about internet regulation, including how states and their societies deal with disinformation? Essentially there are two aspects to this: first, the jurisdictional challenge means that states find it increasingly difficult to effectively apply and

---

<sup>167</sup> See fn 86, 163

<sup>168</sup> See further J. Hörnle *Internet Jurisdiction Law and Practice* (OUP 2021) 4

<sup>169</sup> Eg Universal Declaration of Human Rights 1948, Art 21 (3)

<sup>170</sup> J. Hörnle *Internet Jurisdiction Law and Practice* (OUP 2021) 81-113

<sup>171</sup> *ibid*

<sup>172</sup> J. Hörnle *Internet Jurisdiction Law and Practice* (OUP 2021) 436

<sup>173</sup> See further J. Hörnle *Internet Jurisdiction Law and Practice* (OUP 2021) 7

<sup>174</sup> S Wheatley "Election hacking, the rule of sovereignty, and deductive reasoning in customary international law" (2023) 36 (3) *Leiden Journal of International Law* 675-689, 677

<sup>175</sup> D Herzog *Sovereignty RIP* (Yale University Press 2020)

<sup>176</sup> H Krüger "Of Zombies, Witches and Wizards- Tales of Sovereignty" (2022) 33 (1) *European Journal of International Law* 275-296, 294

<sup>177</sup> Arts 1, 3 of the Montevideo Convention on Rights and Duties of States, 26. December 1933, 165 LNTS 19 and Art 2 (1) of the UN Charter "principle of the sovereign equality of all its Members"; see also H Krüger "Of Zombies, Witches and Wizards- Tales of Sovereignty" (2022) 33 (1) *European Journal of International Law* 275-296, 287

enforce their national laws against disinformation online. This entails that states lose sovereignty over the governance of disinformation. Secondly, where disinformation campaigns are actively launched by foreign states to interfere with elections and political processes, this raises the question whether this is a type of foreign interference, illegal under *international* law. The concept of sovereignty might help us to understand what measures states can take, and to understand global social media companies' power struggle with states. The UK Parliamentary Foreign Affairs Committee is currently holding an Inquiry mapping foreign interference by state and non-state actors, precisely in order to understand how this affects diplomatic relations and the actions a state can take.<sup>178</sup>

#### 4.1 The Jurisdictional Challenge

The ubiquitous and cross-border nature of the internet has led to a massive jurisdictional challenge. In principle, speech and social interactions online are without borders and take place remotely without necessarily being tied to any particular territory. This means that jurisdictional claims overlap, several states may try to apply their national law, in practice it may be difficult to enforce the civil and criminal law across a border, and there is a real challenge in investigating illegal acts online.<sup>179</sup> This means that the power of states (sovereignty) to govern is reduced and conflicts of law abound (jurisdictional challenge to effectively apply and enforce the law). Multiple states may assert jurisdiction, laws differ in substance, there is confusion as to the applicable law(s) with significant overlap and states find it difficult to apply the law reflecting their constitutional values to the activities of their citizens.<sup>180</sup> As a consequence, states exert pressure on online service providers as the gatekeepers of the internet to recreate borders online by identifying the physical location of the users of their online services. Generally speaking, states force online service providers to apply EU law to users located in the EU, UK law to users located in the UK, and US law to users located in the US. This leads to the re-establishment of borders on the internet<sup>181</sup>. I argue, though that this is inevitable in our current legal system where the law is largely that of nation states and where the ubiquity of the internet clashes with a system of regulation and values based on nation states.

States' controlling of speech, such as disinformation online, has largely been presented negatively as extending state repression against free speech into cyberspace.<sup>182</sup> However the recreation of national delineations for determining jurisdiction for the purposes of state regulation over citizens' internet activities and data is neither repressive nor liberal *per se*. It can equally be represented as protecting national values, including protecting Constitutional rights, media pluralism, civil liberties and the rule of law, but it can equally constitute the extension of undemocratic political power and repression of opposition.<sup>183</sup> The extension of state regulation to disinformation can be explained by the concept of sovereignty as a

---

<sup>178</sup> UK Parliament (15. January 2025) <https://committees.parliament.uk/committee/78/foreign-affairs-committee/news/204722/new-inquiry-disinformation-diplomacy-how-malign-actors-are-seeking-to-undermine-democracy/>

<sup>179</sup> See further J. Hörnle *Internet Jurisdiction Law and Practice* (OUP 2021) Chapters 3-6

<sup>180</sup> *ibid*

<sup>181</sup> BA Simmons et al "Cyberborders Exercising State Sovereignty Online" (2023) 95 *Temple Law Review* 617-640, 624-626

<sup>182</sup> BA Simmons et al "Cyberborders Exercising State Sovereignty Online" (2023) 95 *Temple Law Review* 617-640, 624

<sup>183</sup> *Ibid* 625

national<sup>184</sup>, societal preference for domestic regulation over *other* states' conflicting conception of free speech overreaching into the domestic sphere:

"(...) border orientation captures how the state/society governs forces emanating from the rest of the world ("horizontal" relations). As such, it demonstrates a preference for defining and protecting the home environment from the global one (...) Despite the technological challenges of regulating the internet, governments are "fencing" the internet for very traditional reasons: they want to maintain sovereign control over the information environment in their national territory."<sup>185</sup>

While both democratic and undemocratic states exercise sovereignty, the concept itself is speech neutral. Sovereignty is concerned with international relations, *ie* the power relationships between states and non-interference with the domestic matters of a foreign state.

However, it is clear that this reintroduction of borders in cyberspace is only partially effective. If online service providers ignore or partially ignore compliance with national law in respect of the services they provide to users in a jurisdiction, states may find it difficult to effectively enforce regulations against disinformation. While big tech companies initially were motivated by the legitimate business purpose of increasing their advertising income, in the current era we see direct political influence in social media companies and exploiting their technology for political power, with Elon Musk's "X" as the paradigmatic example. Likewise, in his announcement to replace fact-checkers *Zuckerberg* also stated that he would work with President-elect Trump to fight "censorship" in other jurisdictions including Europe and Brazil.<sup>186</sup> Given the ideology and political aims of these platforms, disinformation and the disruption caused is intentional, as they are seeking political power themselves and this course puts them in confrontation with national governments. So while states seek to introduce borders in internet communication this is not always effective, as states jurisdiction is challenged by limited investigative and enforcement jurisdiction. A number of US based platforms have failed to comply with the DSA and the European Commission has started enforcement proceedings<sup>187</sup>- this is a good example of the power struggle between big tech companies and the state. The outcome of this particular struggle is not yet clear.<sup>188</sup> Another example is of course the bipartisan US legislation<sup>189</sup> ordering that Tik-Tok be sold to a US company or banned in the US from 19. January 2025. While this law is based on concerns of data harvesting of US users and national security concerns<sup>190</sup>, not disinformation, this is also an example of recreating borders and exerting data sovereignty,

---

<sup>184</sup> Or in the case of the EU a regional preference

<sup>185</sup> Ibid 625

<sup>186</sup> G Scofield "Alliance between Meta and Trump is likely to create informational, economic and geopolitical conflicts around the world" (8. January 2025) The Conversation <https://theconversation.com/alliance-between-meta-and-trump-is-likely-to-create-informational-economic-and-geopolitical-conflicts-around-the-world-246872>

<sup>187</sup> EU Commission Press Release (12. July 2024) in relation to X [https://ec.europa.eu/commission/presscorner/detail/en/ip\\_24\\_3761](https://ec.europa.eu/commission/presscorner/detail/en/ip_24_3761) ; EU Commission Press Release (16. May 2024) in relation to Meta [https://ec.europa.eu/commission/presscorner/detail/en/ip\\_24\\_2664](https://ec.europa.eu/commission/presscorner/detail/en/ip_24_2664)

<sup>188</sup> The European Commission has extensive investigatory and sanctioning powers, for example fines of up to 6% of the worldwide annual turnover of the company, as well as periodic penalties of up to 5% of the average daily worldwide turnover for each day of delay of compliance, see FN 124

<sup>189</sup> Protecting Americans from Foreign Adversary Controlled Applications Act, April 2024

<sup>190</sup> The US Supreme Court applied intermediate scrutiny on the basis that this was not a law *about* content, see p.5

in order to deal with foreign interference. The legislation has been upheld as compatible with the 1<sup>st</sup> Amendment by the US Supreme Court<sup>191</sup> and at the time of writing Tik Tok was blocked in the US, to the frustration of its 170 million users.

#### 4.2 Foreign Interference as an Act against State Sovereignty?

Disinformation can stem from national bad actors, but equally can be an act of interference by foreign states with the political system and election integrity. Examples are *Elon Musk* interfering with misleading comments on British politics<sup>192</sup> or the election campaign in Germany<sup>193</sup> or Russian interference with elections in Western countries<sup>194</sup>.

This raises the question whether this is illegal under international law as an act against the sovereignty of the states concerned. This is controversial.<sup>195</sup> The Talinn Manual 2.0 provides in Rule 4: "States must not conduct cyber operations that violate the sovereignty of another State."<sup>196</sup> However, Rule 4, and the scope of sovereignty over the internet is unclear and disputed.<sup>197</sup> Some states deny that cyber operations interfere with a state's sovereignty, whereas other states claim that cyber operations can indeed be an interference with a state's sovereignty.<sup>198</sup> Cyber operations in this context mean hacking, computer misuse and foreign interference.

*Wheatley* argues that usurping governmental functions<sup>199</sup> violates the rule of sovereignty<sup>200</sup>, but merely interfering with the exercise of government processes (such as elections) does not.<sup>201</sup> This means that foreign state cyber-operations interfering with the outcome of elections through the spreading of disinformation in the target state are unlikely to be an interference of a state's sovereignty illegal *under international law*. Likewise, a state colluding with a social media platform to exert influence over a population in a foreign state would not count as an act of illegal interference under *international law*.<sup>202</sup> However the *national law* of a state may criminalise such conduct. The UK has created a criminal offence

---

<sup>191</sup> *TikTok v Merrick B Garland* US Supreme Court Westlaw 17. January 2025

<sup>192</sup> The Guardian (20. November 2024) <https://www.theguardian.com/media/2024/nov/20/mps-summon-elon-musk-x-role-uk-summer-riots>

<sup>193</sup> M Fitzpatrick "The far-right is rising at a crucial time in Germany, boosted by Elon Musk" (1. February 2025) The Conversation <https://theconversation.com/the-far-right-is-rising-at-a-crucial-time-in-germany-boosted-by-elon-musk-247895>

<sup>194</sup> BBC News article about RT paying a covert US operation to influence the outcome of the 2024 Presidential elections and criminal prosecutions (4. September 2024) <https://www.bbc.co.uk/news/articles/c8rx28v1vpro>

<sup>195</sup> S Wheatley "Election hacking, the rule of sovereignty, and deductive reasoning in customary international law" (2023) 36 (3) *Leiden Journal of International Law* 675-689, 676

<sup>196</sup> The Tallinn Manual 2.0 (2017) is a work of international law scholarship digesting the principles of international law on cyber-operations which do not use force or armed conflicts, published by Cambridge University Press, see <https://www.cambridge.org/gb/universitypress/subjects/law/humanitarian-law/tallinn-manual-20-international-law-applicable-cyber-operations-2nd-edition?format=PB>

<sup>197</sup> FN 195, 680

<sup>198</sup> FN 195, 682-683

<sup>199</sup> Using deductive reasoning, in the sense of "exclusive right to exercise sovereign authority with respect to a territory" p. 688, referring to *Island of Palmas* (Netherlands v. USA), Award of 4 April 1928, 2 RIAA 829 (1928), at 838

<sup>200</sup> For example, the *Corfu Channel* case, (United Kingdom of Great Britain and Northern Ireland v. Albania), Merits, Judgment of 9. April 1949, [1949] ICJ Rep. 4

<sup>201</sup> See FN

<sup>202</sup> TIK TOK

of foreign interference, where a foreign power spreads disinformation interfering with political processes such as elections in the UK.<sup>203</sup>

Arguably international law *should* be applied against illegal cyber-operations, including foreign interference and state practice should be changed in this respect. Of course it is uncertain whether state practice will ever emerge in this respect. But the concept of state sovereignty might help to explain the authority of the state to contain the power of social media companies and it is to this issue that I will turn next.

### **4.3 State Sovereignty and Social Media Companies**

This article argues that a new powerful type of actor has emerged on the international political scene in addition to states, namely the social media companies themselves. International law may not yet recognise this new category of international actor. But state sovereignty should additionally be concerned with the power relationships between the state and multi-national big tech companies, given their power and influence over the world's citizens. The key question, though, in this power struggle is who is the better guarantor for our civil liberties, including but not limited to free speech: the (democratic) state with a government constrained by a Constitution, the rule of law and mechanisms protecting fundamental rights, or private social media companies having no fundamental rights guarantees and giving unedited, but highly manipulated voices<sup>204</sup> to the "man (and woman<sup>205</sup>) on the street" and increasingly turning into an oligarchy of unelected, but extremely rich and powerful big tech companies?

Civil liberties originally conceived against the all powerful state, are now threatened by entities more powerful than the state and controlled by CEOs with political power aspirations. This threat raises the question of whether the state should be the guarantor of fundamental rights and protect our civil liberties against unaccountable big tech and powerful social media platforms. This is how the rule of law and fundamental rights connect with questions of sovereignty. Democratic states possess institutions and mechanisms to safeguard civil liberties. If these institutions and mechanisms are undermined by the manipulation of big tech companies there is a danger that the democratic state itself will disappear. Social media platforms arguably are more powerful in terms of political influence than any state sovereign and the traditional conception of state sovereignty under international law does not take this into account. Therefore, platform regulation can be seen as a power struggle between nation states and platforms. This power-struggle is a threat to the constitutional order of states, the rule of law, and the guarantee of civil liberties.

The power struggle between big tech companies and the state is challenging democratic nation states' sovereignty over their citizens' activities and data. This power struggle may in the past have been conceived as greater liberty of cyberspace vis-à-vis the state, but citizens' civil liberties are now threatened by the power of big tech controlling their data and

---

<sup>203</sup> Sections 13-15 National Security Act 2023, for or on behalf of a foreign power, s.31 (1) (a)

<sup>204</sup> Manipulated, because they target content to us based on behavioural profiling, and maximising engagement through predicting our subliminal desires

<sup>205</sup> Given the current waves of misogyny and abuse directed at women, arguably women have much less of a voice than men in the social media space, creating new inequalities and serious questions about free speech for women at the expense of that of men. But this is a different topic to be explored in another article.

their online behaviour through opaque and unaccountable algorithms which manipulate the information citizens see.<sup>206</sup>

## **5. Conclusion**

While the ubiquitous nature of the internet may have increased citizens' access to information and freedom of expression generally, this new freedom is increasingly being severely curtailed by a handful of powerful big tech companies. The purpose of content regulation is to ensure the equivalent protections against social media companies. It is in this light that we need to perceive the fight against disinformation. Therefore, state sovereignty should mean a right for the state to regulate social media companies to protect its citizens' right to receive information and to protect media pluralism. Rather than prohibiting "censorship", international law and human rights frameworks need to limit *how* states regulate social media companies and ensure that this regulation enables plurality, review of decisions, complaints mechanisms, transparency, access by researchers and other safeguards.

---

<sup>206</sup> See also the UK Parliament Science and Technology Select Committee's Social Media Inquiry