



# GHOST IN THE MACHINE

Addressing the consumer harms of generative AI

# TABLE OF CONTENTS

<b>EXECUTIVE SUMMARY</b>	<b>5</b>
<b>1 – INTRODUCTION</b>	<b>6</b>
<b>1.1 An overview of generative artificial intelligence</b>	<b>7</b>
1.1.1 Examples of generative AI models	8
1.1.2 The generative artificial intelligence actor chain	11
1.1.3 Open source or closed source models	12
1.1.4 General purpose AI	12
<b>1.2 Consumer applications</b>	<b>13</b>
<b>2 – HARMS AND CHALLENGES OF GENERATIVE ARTIFICIAL INTELLIGENCE</b>	<b>14</b>
<b>2.1 Structural challenges of generative AI</b>	<b>15</b>
2.1.1 Identifying the concrete risks of generative AI	15
2.1.2 Technological solutionism	17
2.1.3 Concentrating power in the hands of big tech	17
2.1.4 Opaque systems and lack of accountability	19
<b>2.2 Manipulation</b>	<b>22</b>
2.2.1 Mistakes and inaccurate output	22
2.2.2 The personification of AI models	23
2.2.3 Deepfakes and disinformation	25
2.2.4 Detecting AI-generated content	27
2.2.5 Generative artificial intelligence in advertising	28
<b>2.3 Bias, discrimination, and content moderation</b>	<b>29</b>
2.3.1 Bias in training data	29
2.3.2 Content moderation	30
<b>2.4 Privacy and data protection</b>	<b>31</b>
2.4.1 Privacy challenges related to data sets used for model training	31
2.4.2 Privacy challenges related to generated content	32

<b>2.5</b>	<b>Security vulnerabilities and fraud</b>	<b>32</b>
<b>2.6</b>	<b>Replacing humans in consumer-facing applications with generative AI, wholly or in part</b>	<b>33</b>
	2.6.1 Challenges related to combining human- and automated decision-making	<b>33</b>
<b>2.7</b>	<b>Environmental impact</b>	<b>34</b>
	2.7.1 Climate impact	<b>34</b>
	2.7.2 Water footprint	<b>36</b>
	2.7.3 Greenwashing & hopes for Green AI	<b>37</b>
<b>2.8</b>	<b>Impact on labour</b>	<b>37</b>
	2.8.1 Labour exploitation and ghost work	<b>37</b>
	2.8.2 Labour automation and threats to jobs	<b>38</b>
<b>2.9</b>	<b>Intellectual property</b>	<b>38</b>
<b>3</b>	<b>REGULATIONS</b>	<b>40</b>
<b>3.1</b>	<b>Data protection law</b>	<b>45</b>
	3.1.1 Data subject rights	<b>46</b>
	3.1.2 The Italian DPA's decision concerning ChatGPT	<b>47</b>
<b>3.2</b>	<b>Consumer law</b>	<b>48</b>
	3.2.1 Consumer law in a U.S. setting	<b>49</b>
<b>3.3</b>	<b>General product safety law</b>	<b>50</b>
	3.3.1 The General Product Safety Directive	<b>50</b>
	3.3.2 The General Product Safety Regulation	<b>50</b>
<b>3.4</b>	<b>Competition law</b>	<b>51</b>
<b>3.5</b>	<b>Content moderation</b>	<b>51</b>
<b>3.6</b>	<b>The draft Artificial Intelligence Act</b>	<b>52</b>
	3.6.1 The EU Commission's proposal	<b>52</b>
	3.6.2 The EU Council's position on the AIA	<b>53</b>
	3.6.3 The EU Parliament's position on the AIA	<b>53</b>
	3.6.4 The Artificial Intelligence Act must protect consumers	<b>54</b>

<b>3.7</b>	<b>Liability</b>	<b>54</b>
	3.7.1 Product Liability Directive	54
	3.7.2 Revised Product Liability Directive	55
	3.7.3 AI Liability Directive	55
<b>3.8</b>	<b>Industry standards and guidelines</b>	<b>56</b>
<b>4</b>	<b>– THE WAY FORWARD</b>	<b>57</b>
<b>4.1</b>	<b>Consumer rights principles that are key for safe and responsible AI</b>	<b>59</b>
<b>4.2</b>	<b>Policy recommendations</b>	<b>60</b>
	4.2.1 Calls to action and empowerment of enforcement agencies	60
	4.2.2 Decision makers – strategic measures	61
	4.2.3 New legislative measures	61
	<b>ENDNOTES</b>	<b>65</b>

---

**Ghost in the machine** – Addressing the  
consumer harms of generative AI

Norwegian Consumer Council  
June 2023

[www.forbrukerradet.no/ai](http://www.forbrukerradet.no/ai)

Design by: Von Kommunikasjon

## Executive summary

There has been an explosion of consumer facing generative artificial intelligence services. These applications can be used to generate synthetic text, images, sound, or video that closely resemble human-made content. As generative AI systems become integrated into popular platforms and tools, adoption of the technology is poised to keep rising. Meanwhile, a number of emerging challenges have spurred numerous debates about how to ensure that generative AI is safe, reliable, and fair.

This report is a contribution to these discussions, and aims to provide policymakers, lawmakers, enforcement agencies, and other relevant entities with a robust starting point to ensure that generative AI does not come at the expense of consumer and human rights. We cannot know for certain how the technology will develop in the months and years to come but believe that the direction of technological advancement should happen on society's terms. Therefore, we present a number of overarching principles that can help define how generative AI systems can be developed and used in a consumer- and human-centred way.

We also strongly urge governments, enforcement agencies and policy makers to act now, using existing laws and frameworks on the identified harms that automated systems already pose today. New frameworks and safeguards should be developed in parallel, but consumers and society cannot wait for years while technologies are being rolled out without appropriate checks and balances.

The first chapter of this report provides an explanation of generative AI and its uses, alongside several examples and illustrations.

In chapter two, we summarize various current and emerging challenges, risks, and harms of generative AI. This includes challenges related to

- Power, transparency, and accountability,
- wrong or inaccurate output,
- using technology to manipulate or mislead consumers,
- bias and discrimination,
- privacy and personal integrity,
- security vulnerabilities,
- automating human tasks,
- environmental impact,
- labour exploitation.

Chapter three contains an overview of the patchwork of existing and upcoming rules and regulations that may apply to the development, deployment, and use of generative AI systems. This is mostly centred around EU legislation, but with some references to ongoing processes in the United States.

The final chapter contains numerous recommendations on how to address the problematic issues of generative AI. This includes:

- Enforcement of existing laws and regulations,
- ensuring sufficient resources for enforcement bodies,
- stronger consumer protection,
- robust government policies,
- new legislative measures,
- strong obligations that cover developers and deployers of generative AI systems.

The image features a minimalist design with several thin, dark blue lines that create a frame around the central text. These lines are composed of vertical, horizontal, and diagonal segments, some ending in small teal dots. The overall aesthetic is clean and modern.

# 1. INTRODUCTION

Consumer facing artificial intelligence systems have been around in various forms for decades, and are used for personalizing social media, filter emails, recommend streaming content, text translation, and much more. Some of these purposes are benign and discreet, and most people may never realize that they are interacting with an AI-powered system.

A new wave of AI-powered systems in consumer facing applications is fast approaching, with the mass deployment and adoption of generative artificial intelligence (‘generative AI’) systems. Generative AI is a subset of artificial intelligence that can generate synthetic content such as text, images, audio, or video, which can closely resemble human-created content. Such systems are poised to change many of the interfaces and content consumers meet today.

In November 2022, a prototype of the chatbot ChatGPT was released to the public. The application quickly gained worldwide attention, becoming the fastest growing digital service of all time within a month of its release.<sup>1</sup> In the following months, other services for generating text, images, sound, and video, were quickly deployed and iterated upon, sparking a sort of arms race for generative AI systems. Consumers were provided access to these content generators directly in web interfaces, while companies started to embed the content generators in their applications and services.

The sudden and widespread deployment and adoption of generative AI systems sparked public discourse about the promises and perils of the technology. The debate has ranged from how generative AI may be used to increase efficiency in the workforce and to spark cre-

ativity, to how it can be used to spread disinformation, manipulate individuals and society, displace jobs, and challenge artists’ copyright.

The discussion about how to control or regulate these systems is ongoing, with policymakers across the world trying to engage with the promises and challenges of generative AI. This report is a contribution to these discussions by providing an analysis of the most pressing issues from the consumer angle, along with a number of possible solutions and ways to proceed from both a legal, ethical, and political perspective. Although we do not pretend to have the answers to all the questions raised by generative AI, we believe that many of the emerging or ongoing issues can be addressed through a combination of regulation, enforcement, and concrete policies designed to steer the technology in a consumer- and human-friendly direction.

As the development of generative AI seems to move at a breakneck pace, the descriptions throughout this report must be seen as a snapshot of an emerging technology. The report was written between February and May 2023, and does not include any new information from papers published after June 1st.

The Norwegian Consumer Council is a publicly funded, independent consumer organization, that represents consumer interests. This report was written with contributions from BEUC, Miika Blinn from VZBV, Kris Shrishak from the Irish Council for Civil Liberties, Daniel Leufer from Access Now, Jon Worth, Marija Slavkovik, and Anja Salzmann from the University of Bergen.

## 1.1 An overview of generative artificial intelligence ---

Generative AI is a blanket term used to describe algorithmic models that are trained to generate new data, such as text, images, and sound. While these applications rely on different types of input data, the general principles behind how they are trained are similar. The emergence of advanced generative AI is possible due to an enormous amount of content available on the internet, combined with advances in machine learning and computing power.

Generative AI models work by analysing large amounts of information to predict and generate the next word in a sentence, feature of an image, etc. This is done by detecting patterns in and relationships between data points in the training data, which in turn allows the system to replicate similar patterns to generate synthetic content, for example a piece of writing, music, or a video clip. This process can also be described as a complex ‘mash-up’ of content from the data the system was trained on. In other words, they are predictive mod-

els that are trained to “connect the dots” between data points in existing content to generate synthetic content.

Generated content is probabilistically and randomly generated based on certain input (or ‘prompts’), which are usually written by a human. Therefore, the output of any given generative AI model is likely to be different for each person prompting the model and may both resemble patterns in the training data or appear to be something completely new.

### 1.1.1 EXAMPLES OF GENERATIVE AI MODELS

There are various types of generative AI models, including large language models (LLM) that can respond to text by generating new text, and multi-modal models that can generate more than one type of output or respond to more than one type of input, for example chatbots that can also generate images when prompted to do so. A short introduction of the most popular generative AI models on the market today is presented below, accompanied by some relevant examples.

#### 1.1.1.1 Text generators

Text generators are a type of generative AI model that can generate text passages based on predictive analysis, which are built on large language models.<sup>2</sup> These models are usually trained on enormous amount of text scraped from the internet, including books, forums, news sites, social media, etc. Text generators can be used for writing essays, coding, chatbots, and augmenting search engines, amongst other things. In many cases, text generators are meant to generate text that appears to be written by a human, for example by generating text written in a first-person perspective, using emojis, or by writing text indicating that it has the capability to experience human emotions. Some text generators are multimodal and can generate text based on images.

Although text generators have existed in some form for several years, for example as prediction tools for typing text messages, the discussion around the technology

gained momentum during the autumn of 2022, with the public release of ChatGPT, which is owned and operated by the company OpenAI (which is also the owner of DALL-E, see below). ChatGPT3 is available online for those who create a free account, while the more powerful ChatGPT4 model is available for a monthly subscription fee.<sup>3</sup>

In January 2023 Microsoft announced a major investment in ChatGPT and launched new features in the Bing search engine powered by the technology.<sup>4</sup> Microsoft has announced that it intends to integrate ChatGPT into its other services, including the Microsoft Office suite of applications, for example to automatically take notes during meetings in Microsoft Teams.<sup>5</sup>

Google has also developed a LLM that can generate text, called LaMDA. In the wake of Microsoft’s investment in ChatGPT, Google rolled out similar features to its search engine with a text generator called Bard.<sup>6</sup> Google is also planning to introduce various AI-powered features such as drafting and summarizing emails, as well as brainstorming and writing documents in its Workplace applications.<sup>7</sup>

AI

Write a one paragraph poem about the consumer issues of generative AI



In the realm where codes unfurl,  
Generative AI, a modern swirl.  
Amidst its marvels and digital grace,  
Lies a quandary, a consumer's chase.  
Ownership blurred, rights cast astray,  
As creations emerge, who holds the sway?  
Questions abound, ethereal hue,  
In this labyrinth of creations anew.

*Poem about the consumer issues of generative AI, ChatGPT.*

<sup>2</sup> Large language models are sophisticated AI models that are designed to generate text that resembles human language. They are normally trained on vast amounts of text sources to “learn” patterns and grammar. LLMs

can be used for tasks such as machine translation, sentiment analysis, human-machine interaction, proofreading, and many other purposes.



Meta has developed the LLM Galactica trained on scientific articles and materials, which is meant to “store, combine and reason about scientific knowledge”. After the model was released as a public demo in November 2022, the public-facing demo was quickly removed due to generating text containing multiple errors and biases.<sup>9</sup> In February 2023, Meta released another LLM, called LLaMa (Large Language Model Meta AI). LLaMa is an open source model, which was initially released to researchers based on an access application process. In March 2023 the model was leaked on a public message board, meaning that anyone with a relatively powerful computer can download, use, and adapt the model.<sup>9</sup>

There are also several open source LLMs that are developed and maintained by smaller actors. For example, the text generator BLOOM is available through the company Hugging Face,<sup>10</sup> while StabilityAI has released open models under the moniker StableLM.<sup>11</sup>

### 1.1.1.2 Image generators

Generative AI models that are trained to generate images can be collectively classified as image generators.

They can create images from text prompts (‘text to image’) or from existing images (‘image to image’). Image generators work by analysing huge amounts of existing pictures, such as photographs, paintings, etc., which are often scraped from various sources on the internet. By training the algorithm on these data sets, the model can generate images of different objects (‘a chair’, ‘a train’), people (‘a young woman’, ‘Jerry Seinfeld’), and styles (‘impressionism’, ‘in the style of Edward Munch’). The most widely used image generators as of June 2023 are Midjourney,<sup>12</sup> DALL-E,<sup>13</sup> and Stable Diffusion.<sup>14</sup>

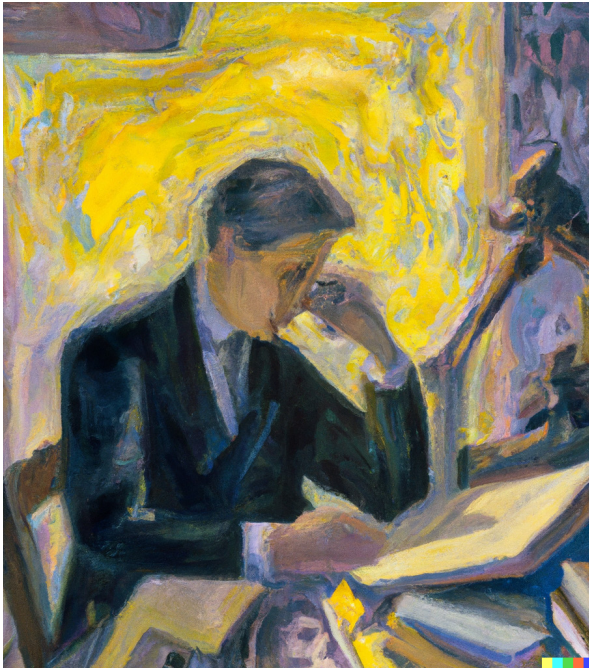
*Hyper realistic photo of political advisor writing a paper on generative ai, in an open office plan, Midjourney.*

Midjourney is available through the chat service Discord. It is possible to join the official Midjourney Discord server, to ask a Midjourney bot to “imagine” a picture based on various prompts. For example, the person prompting the system could type “/imagine hyper realistic photo of political advisor writing a paper on generative ai, in an open office plan”. The bot responds in the chat with four generated pictures. While Midjourney was free to try for the first few months after its release, for a limited number of generated images, it has since become a paid subscription service.

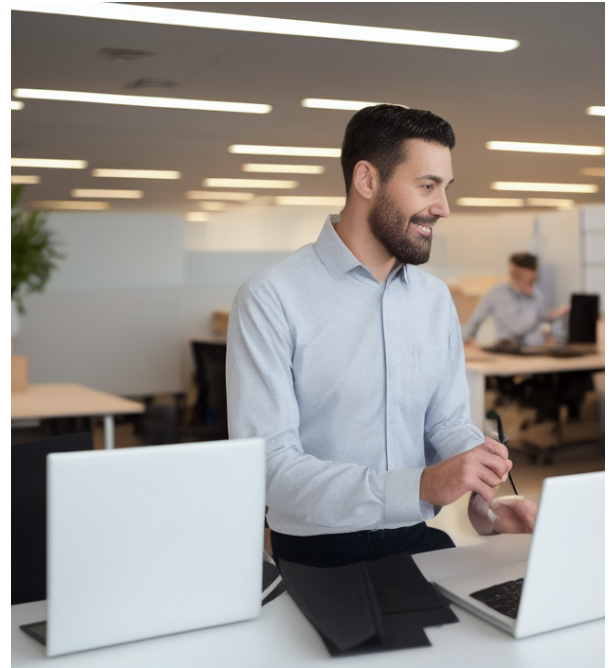
The company Midjourney Inc. owns the generative AI model that generates the images and runs and controls both the model itself and the servers that it is hosted on. This means that the company can restrict access, change the model, and add content filters to control what kind of images the model can and cannot generate.

The generative AI model DALL-E is accessible through the website of its owner OpenAI. Individuals can create an account and receive a limited number of tokens each month that can be used to generate images. Images are





*An oil painting by Munch of a policy advisor writing a paper on generative AI, DALL-E.*



*A photo of a consumer policy advisor writing a paper about the consumer challenges of generative AI, in an open office plan, Stable Diffusion 1.5.*

generated by entering different prompts into the website interface. If the person runs out of tokens, they can pay to receive more. Like Midjourney, the model behind DALL-E is owned, controlled, and operated by its parent company.

The generative AI model Stable Diffusion is developed by the company Stability AI. Unlike Midjourney and DALL-E, Stable Diffusion is an open source model that can be freely downloaded by anyone, and does not require a subscription or access to the internet to use. Once the necessary software is downloaded, unlimited images can be generated locally. Running Stable Diffusion locally only requires a computer with relatively powerful consumer-grade graphics card.

Stability AI only trains and distributes the basic models for Stable Diffusion, while it is possible for anyone with access to them to continue training and developing new models that are based on the original Stable Diffusion models. These new models can then be distributed to others. In practice, this means that the company does not control the model nor its output.

#### 1.1.1.3 Audio generators

Audio generators use generative AI technologies to create audio clips based on text prompts (for example

text to voice). Such models are trained on existing voice data, music, etc. Audio generators can be used to create AI-generated music<sup>15</sup> and voices, and there are models capable of recreating the voice and pitch of individual individuals.<sup>16</sup>

For example, the company ElevenLabs has released a model that allows anyone to convert short text input into voice clips, in a selection of different voices.<sup>17</sup> Microsoft has announced the generative AI model VALL-E, which the company claims can generate realistic voices based on a three second voice sample.<sup>18</sup> As of May 2023, VALL-E has not been released to the public.

#### 1.1.1.4 Video generators

Video generators can be used to create video clips based on text prompts (text to video), images (image to video), or existing clips (video to video). As it is more complicated to generate authentic looking video footage than to make still images, this technology is somewhat less developed at the time of writing.

However, this may change in the near future, as several major companies are actively working on models for generating video.



Frame from video generated in Make-A-Video, Meta AI.<sup>19</sup>

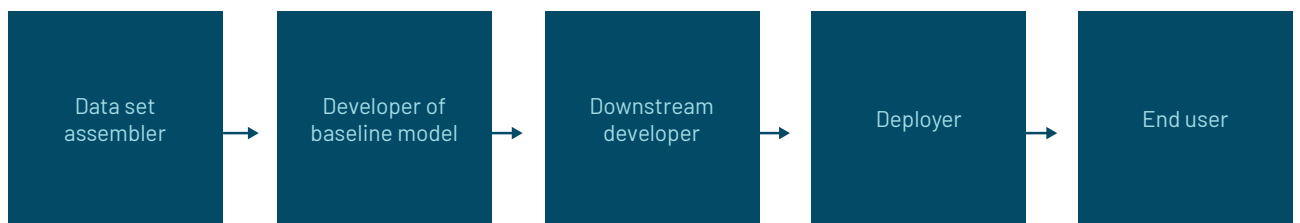
Meta has developed a model to turn short texts into video clips,<sup>20</sup> and Google has announced a similar system.<sup>21</sup> As of May 2023, neither of these systems have been made available to the public.

Stability AI, the company behind Stable Diffusion, has released a model for generating animations from text prompts and images.<sup>22</sup> The company Runway has released a mobile app that can be used to generate short video clips from existing videos.<sup>23</sup>

### 1.1.2 THE GENERATIVE ARTIFICIAL INTELLIGENCE ACTOR CHAIN

From assembling data sets for training models to deploying and prompting generative AI systems, there are potentially many different actors. These actors may all influence the system or how it is used in various ways. The relevant actors are shown in the illustration just below.

Each box represents a different actor in the actor chain. These actors may all be within one organisation but will often be spread across several organisations. The arrows are only pointing one way to maintain a simple representation; it is of course possible that there could be feedback loops between different actors.



The data set assembler collects and systemizes the necessary data to train a generative AI model. There are numerous available open source data sets that have been compiled and labelled by scraping a huge number of sources online. In many cases, such data sets are compiled for research purposes, and made freely available. A company developing a generative AI model can thus train their model on datasets that someone else has assembled. The developers of generative AI models create a baseline model, which may then be trained and tuned for certain more specific contexts or applications by downstream developers. In some cases, this fine-tuning of the model may be done by the same actor as the training of the baseline model, while in other cases, the fine-tuning may be done by a separate actor entirely. This may be another company, or in the case of open source models, the model can be fine-tuned by anyone with a relatively powerful computer.

To further complicate the matter, as general-purpose generative AI models are being integrated into other applications, the company or entity deploying the system may be separate from the company developing the model and/or fine-tuning it.

Finally, there are end users, engaging with the deployed model. In consumer use cases, the consumer will typically be an actor, by prompting the model to generate, for example, a text or an image. Consumers may also indirectly encounter generative AI systems when interacting with a business, for example if a customer service agent uses a text generator to generate answers to consumer queries, or if a consumer is prompted to make certain queries based on AI-generated recommendations.

There are numerous actors in the generative AI system actor chain. It is crucial to understand the relationship between these actors to understand how generative AI should be regulated, and at which point the actor chain different harms arise.

Illustration of different actors in the generative AI actor chain.

### 1.1.3 OPEN SOURCE OR CLOSED SOURCE MODELS

There are significant differences in how generative AI models are distributed and controlled. Many models are proprietary closed source models that run on cloud servers controlled by the owner of the system (the “system owner”). This means that consumers can access the model through the internet, and that the system provider can change the model at any time, add content filters, restrict access, etc. For closed systems, the system owner provides the processing power required both to train the model, and to generate synthetic content.

For a closed source generative AI model, it is not possible to know how the model works, what data it was trained on, and how parameters are weighed, unless the company behind the model publishes sufficient documentation or provides access to third party auditors, enforcement agencies, or researchers. In many cases, such information may be kept secret due to security and/or business interests.

On the other hand, some generative AI models are released as open source, which may take various forms. Different parts of the system, such as the data set, source code, model parameters and weights, may be made available to third parties.

When the source code is made available to the public, it can be used, studied, tested, modified, and distributed by anyone. This means that it can be inspected for errors and vulnerabilities. It can also be improved and iterated upon collaboratively. Using, modifying, and distributing open source software is generally governed by license terms. Given the open source software and the data set, anyone with sufficient computing resource could reproduce the generative AI model, although the computing resources needed to do so are significant enough that in practice it will likely be limited to large companies.

However, for generative AI systems to be truly open source, the model itself should be made available to the public. In such a case, applications developed based on these models might also be released as open source application, such as the image generator Stable Diffusion, or it can be adapted into a closed source application.

Open source generative AI models can be downloaded by anyone. With a powerful enough computer, it can be used to generate data and update the model as one pleases. The source code, parameters, etc. of such models can be

inspected by anyone. However, this does not necessarily mean that they can understand how the model works, as sufficiently complex models may be more or less impossible to understand the inner workings of.

Once an open source generative AI model is released to the public, there is practically nothing the developer of the baseline model can do to influence how the model functions. This means that any content filters and other artificial limiters placed on the model may be altered or removed by downstream developers or deployers. This creates both benefits and significant drawbacks, which will be elaborated on below.

### 1.1.4 GENERAL PURPOSE AI

While some AI models are designed with a specific purpose and use case in mind, such as early discovery of cancer cells, many generative AI models are examples of so-called ‘general purpose artificial intelligence’. This means that the basic system, such as a large language model, is trained to be able to respond to a vast variety of situations and interactions and can be adapted to be used in new contexts.

Unlike a model with a specific purpose, it is extremely difficult or impossible for the developers of a general purpose AI model to foresee the possible uses and abuses of the technology. This makes it particularly important that such models are subject to technical, scientific, legislative, and regulatory scrutiny before they are widely adopted. However, applications such as ChatGPT have already been released to the wider public without rigorous evaluation, impact assessment or scrutiny, while being increasingly opaque and inaccessible to third party auditors and researchers.<sup>24</sup> It is worth considering whether this leads to a desirable future, considering the many harms touched upon in chapter 2 of this report.

## 1.2 Consumer applications

Several different types of generative artificial intelligence are publicly available to consumers today. Many of these are readily available to use by anyone with an internet connection, and do not require expert technical knowledge to use. Some of them are directly accessible through web interfaces, while generative AI technology is also increasingly being integrated into other services such as online search, learning and administration software, and social media.

As of May 2023, the most popular consumer uses for generative AI models are text and image generation. However, with major consumer-facing companies such as Microsoft, Meta, and Google investing heavily in the technology, the use cases are likely to expand in the coming months, as generative AI models are implemented in various services.

For example, text generators can be a useful tool to streamline and/or optimize mundane tasks, functioning as a kind of multi-purpose digital assistant. This may include changing internet search functions, automating certain tasks such as writing code, transcribing voice messages, or personalizing services in various ways. While such applications may be useful and efficient in certain contexts, there are also significant risks and drawbacks, which will be explored further in the following chapters.

As the technology is developed and adopted, generative AI may be used to automate tedious and time-consuming processes that previously had to be done manually, for example by writing concise texts, filling in forms, generating schedules or plans, or writing software code. It has the potential to make services more cost-efficient, which may also lower costs for consumers, for example when soliciting legal advice.<sup>25</sup> On the other hand, the proliferation of low-cost AI-generated content may replace human labour and human-generated content, thus lowering the quality of consumer-facing services such as customer support. The technology also opens new avenues for consumer manipulation in areas such as advertising or product recommendations and can facilitate or obfuscate discriminatory practices.



# 2. HARMS AND CHALLENGES OF GENERATIVE ARTIFICIAL INTELLIGENCE

**“From violations of privacy and personal integrity to the creation of fraud and misinformation, generative AI models introduce vast risks and challenges, while turbocharging others.”**

There have been several controversies surrounding the development and use of generative artificial intelligence. From violations of privacy and personal integrity to the creation of fraud and misinformation, generative AI models introduce vast risks and challenges, while turbocharging others. Concrete and highly relevant examples of this are chatbots and search engines providing incorrect but convincing information, the abuse of cheap labour in the global south for content moderation, and a significant environmental impact due to resource consumption. It is essential that these problems are sufficiently addressed by enforcing, applying, and establishing laws and regulations that serve to protect consumers from various negative consequences.

The issues discussed throughout this report are often not new or unique to generative AI. Algorithmic computer systems have existed for a century, while the technology popularly referred to as artificial intelligence has been around since the 1950s. In the 1960s, the computer scientist Joseph Weizenbaum created ELIZA, a model

that simulated human interaction, using rule-based algorithms.<sup>26</sup> People who interacted with ELIZA attributed human attributes and emotions to the model, even though they were informed that the system had no such capability, mirroring some use cases of generative AI-powered chatbots today.

Issues related to content moderation, algorithmic bias, privacy, and disinformation have been debated at almost every junction as digital technology evolves and is widely used. However, the deployment and public adoption of systems such as ChatGPT, both for technically adept consumers and the general public, alongside its ease of use and wide-scale availability, means that many of these issues have become urgently relevant to analyse from a consumer perspective. As described in the following chapters, a number of these issues may be addressable by enforcing applicable laws and regulations, while other may require other solutions or remedies.

## 2.1 Structural challenges of generative AI

At a basic level, generative AI models are fundamentally designed to reproduce existing material, although in potentially novel ways. This means that such models are inherently prone to reproducing existing biases and power structures. Therefore, while the models have no understanding, mind, or intention of their own, the decision to develop, deploy and use them is inherently political. It is not sufficient to ascribe neutrality or objectivity to the operations or outputs of a generative AI model, because its training data and algorithms stems from human beings, with all that this entails.

As generative AI models are being introduced into all sectors of society, so far with little or no regulatory oversight, there are fundamental issues that need to be addressed. Generative AI models are dependent on large amounts of data that is taken from a multitude of sources, usually without the knowledge or consent of the originator of the data, be it a piece of art, a news article, or a selfie. Information is siphoned and gathered to be used in different ways, with an end goal of enriching a small number of companies. This raises questions of value distribution, usage permission, privacy, accountability, intellectual property, and human rights.<sup>27</sup>

### 2.1.1 IDENTIFYING THE CONCRETE RISKS OF GENERATIVE AI

As with any new technologies, the discourse around generative AI is muddled with a mix of facts, concerns, and a lot of hype and enthusiasm.<sup>28</sup> Many AI systems are being touted as being capable of solving almost any task, often without evidence to back up the claims, a phenomenon that can be described as 'AI snake oil'.<sup>29</sup> When addressing problematic and hazardous issues with the technology, it is important to be able to sort facts from fiction.

**“At a basic level, generative AI models are fundamentally designed to reproduce existing material, although in potentially novel ways. This means that such models are inherently prone to reproducing existing biases and power structures.”**

**“There are serious concerns that focusing on potential long-term scenarios are drawing attention away from many of the current pressing issues of generative AI, potentially leaving these issues insufficiently regulated.”**

TERM	DEFINITION
<b>GENERATIVE ARTIFICIAL INTELLIGENCE</b>	Generates synthetic content based on patterns and structures learned from training data. Used to generate text, images, audio, and video.
<b>GENERAL PURPOSE ARTIFICIAL INTELLIGENCE</b>	Umbrella term for AI systems designed to perform a wide range of tasks across different domains.
<b>ARTIFICIAL GENERAL INTELLIGENCE</b>	Hypothetical AI system that demonstrates human-level intelligence and autonomy. Does not currently exist.

Heavily publicized warnings about the dangers of artificial intelligence have concentrated on hypothetical risks of developing an artificial general intelligence (AGI), meaning a system that is able to perform intellectual tasks that are comparable to the ability of human beings. Theoretically, such systems should be able to think and reason, and be able to perform a broad range of tasks that equal human capacity for thinking. This differs considerably from generative AI models, which have no such capacities. As AGI systems do not currently exist, and there are serious disputes about whether they can ever be realized, such systems will not be further considered in this report.

There have been calls for voluntary moratoriums or ‘pausing’ of developing generative AI models. Some of these calls have focused on a potential future where AI systems have become so powerful that they pose an existential threat to humanity.<sup>30</sup> While such calls acknowl-

edge issues related to accountability, safety, and control over AI systems in general, there are serious concerns that focusing on potential long-term scenarios are drawing attention away from many of the current pressing issues of generative AI, potentially leaving these issues insufficiently regulated.<sup>31</sup>

The argument that a hypothetical general artificial intelligence is an existential threat to humanity implies that concerns about current issues such as discrimination, privacy, and fairness, are inconsequential and marginal.<sup>32</sup> In other words, the narrative concerning a potential “AI supermind” may serve as a distraction from pressing issues that are already present in today’s application of generative artificial intelligence. It is crucial that narratives about hypothetical existential threats to humanity do not come in the way of proposing concrete solutions to the very real issues posed by the technology that exists today.<sup>33</sup>



### 2.1.2 TECHNOLOGICAL SOLUTIONISM

Artificial intelligence is often lauded as the solution for a vast number of issues across sectors, from healthcare and public administration to legal assistance. While this narrative is attractive to both private enterprises looking to sell software solutions, and to policymakers searching for simple remedies to political or regulatory ailments, it needs critical examination.

The belief that almost any issue can be improved or solved using technology is known as 'technological solutionism'. A term coined by the tech critic Evgeny Morozov, technological solutionists tend to gloss over complex and multifaceted social problems in favour of simple mathematical or engineering solutions.<sup>34</sup> This reductionist belief is attractive to service-providers because it allows them to advertise miracle cures – AI snake oil – and to policymakers because technological quick fixes are tangible and usually appear more cost-efficient than examining complicated and often deep-rooted social and political conflicts or inequalities.

As Morozov points out, technological solutionism is dangerous because it often simply does not work. By presenting multifaceted and complex issues as a mere engineering issue to be solved in a lab, solutionism misrepresents social problems and misses the underlying causes. When presented as problems that can be solved by technology, solutionists tend to disregard the social, political, and cultural context that is the backbone of our societies.

When considering the proliferation of artificial intelligence models that are rapidly being deployed across sectors, it is worth keeping the folly of technological solutionism in mind. This is particularly important if generative AI models or similar technologies are being pushed as a remedy or solution to inequalities, such as providing access to mental health tools to people who otherwise would not be able to afford it. While it may appear attractive to outsource mental health care to a large language model that can be deployed and accessed at a relatively low monetary cost, this approach risks reducing the complexity of mental health and the value of human contact to a question of predictive analysis and language modelling.<sup>35</sup> Similarly, before deciding to deploy a text generator as a solution for overworked case handlers in the public sector, for example, it is crucial to consider the context and causes of the problem, rather than adopting developing technologies as a blanket solution.

If such technological quick fixes are adopted at the expense of investment into proven effective measures, which are often costly and difficult to implement, this may come at significant costs to marginalized groups who risk being deprived of effective treatment and measures because their issues are purportedly being addressed by the use of technology.

### 2.1.3 CONCENTRATING POWER IN THE HANDS OF BIG TECH

At the base of the discourse around generative artificial intelligence is a question of power. Generative AI models are products of cultural and political contexts, a context that is embedded in everything from the decision to develop the model, the choice of training data, the tuning of models, and the given purposes for deployment. As such, the already powerful can potentially entrench existing power structures through the technology, while the disenfranchised will remain so unless there is outside intervention. This becomes apparent when a generative AI model generates biased or discriminatory content, but also manifests in aspects such as content moderation practices and in who has access to the systems.

As generative AI models are often trained on data collected from any available sources, some actors are raising questions about whether private companies should be allowed to use the collective knowledge of humanity to turn a profit. The vast amount of information that can be found openly available online can be described as a 'digital commons', as it is a body of resources where practically everyone is a contributor, from individual pieces of data to the public infrastructure of the internet. If the digital commons are siphoned to develop and train proprietary models, this raises ethical concerns about how value generated on the basis of these common resources should be distributed.<sup>36</sup> These concerns extend to data governance issues regarding who should control how data is used, such as if a tech company wants to commercialize AI models trained on indigenous languages.<sup>37</sup>

The question of who controls the development and training of generative AI models and how they are used is of fundamental importance. Those who control the technology have significant potential to create dependencies, set the terms of use, and decide who has access. This entrenchment of power creates overarching concerns about leading tech companies becoming gatekeepers that can exclude rivals and otherwise abuse their increasingly dominant market positions.<sup>38</sup> While

open source models may lower the barrier of entry for certain types of generative AI,<sup>39</sup> such models are still largely dependent on a foundational model that has been developed by actors with access to significant computing resources and training data.

This means that already dominant technology companies such as Microsoft, Google, and Meta, are well positioned to seize the market for generative AIs. With proprietary closed model, the system owner has control over who can access the technology, what it costs, its features, and how it may be used. This may in turn affect academia, attracting what could have been independent researchers to work within the closed domains of big technology companies, where they have access to state-of-the-art technology in the field. Overall, this means that the tech giants can further leverage their dominant positions across different online markets in the field of generative AI.<sup>40</sup>

With only a few generative AI models available on the market, these models are integrated into a variety of services, providing the model owners with significant power. Models can be patched or otherwise modified, functionality can be added or removed, and content can be banned, filtered, or otherwise restricted. If the system owner sets the terms for how its technology may or may not be used, end users or third-party companies integrating the model are at the mercy of the owner.

Some competition concerns may be somewhat alleviated in the case of open source generative AI models, which are not necessarily beholden to the business model, objectives, or whims of a model's original creator. However, even in such cases, many companies do not have the means to compete with big tech firms in offering generative AI solutions to consumers. Large companies benefit significantly from network effects, as more users means more data, which leads to better services. In cases where models are further trained on consumer interactions or feedback, the companies can further improve and fine-tune the models at a rate that is unattainable for smaller competitors.

Dominant actors can further entrench their power by integrating generative AI into their own services that are

already used by millions across the world. For example, by rolling out its chatbot Bard as part of its search engine, Google already has a massive global user base that can be leveraged to boost the adoption of the chatbot. Similarly, as Microsoft implements ChatGPT-based models into its Office suite of applications, the company already has a user base that competitors can only dream of.

Companies can also leverage their market position by making the usage of a generative AI model dependent on using a different service from the same company, by 'bundling' services together. For example, to have access to Microsoft's Bing chatbot functionality, consumers must use Microsoft's Edge browser.<sup>41</sup>

The integration of generative AI models into services such as search engines can also significantly limit consumer choice. For example, in a regular online search

engine, the consumer is presented with numerous search results that they may choose between. If the search engine is replaced by a text generator that provides a single answer to any query, this potentially limits the information available. If similar models

are used for online shopping, this creates new avenues for platforms to self-preference products, by ensuring that the platform's preferred product is the only or the primarily suggested purchase. If the consumer queries "what is the best coffee machine for to suit my needs?", it will be necessary to monitor and control how a chatbot or "shopping assistant" lands at a particular result or recommendation.

### 2.1.3.1 *Walled gardens and its downstream effects*

In order to maximize consumer engagement, many digital service providers have a financial incentive to keep consumers on their platforms as long as possible. This goal can be attained by integrating and bundling as many services into the platform as possible, while creating barriers such as not providing service interoperability to disincentivize consumers to leave the platform. Platforms and services designed to keep the consumer from leaving are called 'walled gardens'.<sup>42</sup>

The integration of generative artificial intelligence into various platforms already seems to facilitate a walled

**“Those who control the technology have significant potential to create dependencies, set the terms of use, and decide who has access.”**

garden approach that may have serious anticompetitive effects on both direct competitors to the large online platforms, and across markets.

Snapchat is for instance introducing recommendations for restaurants or recipes in its AI-chatbot,<sup>43</sup> which would reduce consumers' need to access other services, such as traditional search engines, for these kinds of queries. This is likely a sign of things to come, as major platforms are slated to integrate generative AI models into their services. As companies compete to develop "killer apps", services that integrate as many functions and purposes as possible, these problems are set to exacerbate. Newcomers may find it increasingly difficult to provide stand-alone services to consumers, as they have fewer reasons to exit their applications. This would serve to concentrate power with already established actors, effectively harming the consumer market.

The integration of generative AI into search engines is sparking major concerns for publishers and advertisers, since such integration can effectively create walled gardens.<sup>44</sup> With traditional search engines, the consumer could search for a topic and be presented with a list of links to websites that contain information about the topic. The consumer will then click on one or more of the links and be redirected to a website. Consequently, the website owner generates revenue by displaying ads to the consumer.

This dynamic may change with the introduction of generative AI. Google is for example experimenting with introducing generated, summarized content in its search engine, which fills up the first page on a smaller screen (such as a phone).<sup>45</sup> In short, if consumers can simply ask a chatbot about a subject and receive answers in the same interface, the incentive to visit a third-party website is reduced. If the consumer does not visit the third-party website, the originator of the content cannot monetize the content by displaying ads.

A lack of traffic poses a problem to publishers, whose content may be scraped to display information in through the chatbot interface, but who may struggle to monetize the content. This can therefore have downstream effects by reducing the incentive to produce quality content, potentially leading to content-production becoming automa-

tized as parts of cost-cutting measures. It is also well-established that power concentration among a few actors is seldomly beneficial to healthy consumers markets.

### 2.1.3.2 Data colonialism

If generative AI models are trained on data sets that were indiscriminately scraped from the internet (digital commons), this may also entail large amounts of data from indigenous and other minority groups. The information can then be repackaged and used in new ways, for example to sell technology or services based on the data back to the groups from whom the data originated. The process by which organizations and corporations claim ownership over data produced or harvested from people is called 'data colonialism'. The concept of data colonialism is highly relevant when discussing the operations of generative AI models.

For example, indigenous communities in New Zealand have expressed concerns regarding the development of large language models being trained on hundreds of hours of Māori indigenous language. Community leaders and researchers fear that "if Indigenous peoples don't have sovereignty of their own data, they will simply be re-colonized in this information society".<sup>46</sup>

The language harvested without consent can be distorted and lead to abuse and deprive communities of their rights. According to indigenous communities, it is not for Big Tech to play with such heritage.

### 2.1.4 OPAQUE SYSTEMS AND LACK OF ACCOUNTABILITY

Models such as large language models are generally very technologically complex, but they are not impossible to understand or explain. There are fundamental scientific principles relating to transparency, peer review and rigorous quality control that apply in fields such as the pharmaceutical and aviation industries, which should also apply to developers of AI models. Information about how training data is collected, how the data is labelled, how testing is performed, what decisions are made regarding content moderation, and the environmental and social impacts of the models are just a few areas where transparency is necessary to ensure that risks are mitigated and that claims about the technology are accurate.

<sup>43</sup> Hypothetical AI system that demonstrates human-level intelligence and autonomy. Does not currently exist.

**“There are fundamental scientific principles relating to transparency, peer review and rigorous quality control that apply in fields such as the pharmaceutical and aviation industries, which should also apply to developers of AI models.”**

third-party access to the source code of AI.<sup>54</sup> Closed-door negotiations affecting lawmakers’ ability to create a well-balanced and consumer-friendly market are highly problematic. This bars civil society and other stakeholders from providing important input and appears at odds with crucial democratic principles.

#### 2.1.4.1 Opaque systems reduce accountability

Unfortunately, there are already tendencies from certain AI developers to close off their systems from external scrutiny. For example, Google has committed to a policy change where the company will only share papers after their research has been turned into products.<sup>47</sup> Microsoft researchers have made grand claims about its own AI systems show signs of artificial general intelligence,<sup>48</sup> while not providing researchers with access to the model in order to verify or dispute the claims.<sup>49</sup> Finally, ChatGPT owner OpenAI has claimed that the company’s AI systems, including what training data is used, how the model works, etc., should not be open for external review because giving others access would pose a competition and safety risk.<sup>50</sup>

While the lack of transparency is an issue that appears across the spectrum of the software industry, OpenAI’s own description of the risks of their products tend to border on existential, with its CEO Sam Altman stating that the company is “scared” by the potential harms that may stem from its own systems.<sup>51</sup> Such claims as a pretence for closing down generative AI systems for external auditing and review are worrying tendencies that can mask a number of downstream effects, and which pose enforcement agencies and researchers major challenges.

Researchers at Princeton University have claimed that OpenAI might be misrepresenting the capabilities of their systems, but this is impossible to prove due to the system being closed to external scrutiny.<sup>52</sup> The researchers warn that this also significantly hinders attempts at reproducibility of any claims made by the company.<sup>53</sup>

#### 2.1.4.2 Trade agreements as barriers to transparency

While companies themselves are attempting to close off their systems from external scrutiny, lawmakers may be increasingly limited from requiring transparency by trade agreements. Internal documents from the EU Commission show that digital trade agreements between the EU and US limit European lawmakers’ ability to require

#### 2.1.4.3 Actor chain transparency

The lack of transparency also becomes problematic when service providers implement third-party generative AI models into their services. This may increase the risk of errors or unexpected behaviour from the model.<sup>55</sup> The developer of the baseline model does not necessarily see or understand the downstream contexts in which the model is used, while the service provider or other downstream developers do not sufficiently understand the limitations of the model.

If the service provider is not privy to the data sets used to train the models, or to how the model actually works, the service provider will not be able to give the consumer an explanation about why a certain output was generated. As supply chains for generative AI systems may be complex, with one actor collecting and labelling data sets, while it can be other actors developing the algorithms, training the model, or integrating it into services, it becomes difficult to attribute liability and accountability to the right entity. For the consumer, this may have a negative impact on the right to an explanation, as well as contestability and general transparency obligations.

For example, the retail bank, payments, and shopping service Klarna has announced a collaboration with OpenAI, with plans to integrate ChatGPT into its services to provide a “highly personalized and intuitive shopping experience by providing curated recommendations”.<sup>56</sup> If this system provides flawed recommendations for products, or ranks products in a skewed way, it will be essential that consumers, not to mention enforcement agencies, are able to access and assess data on how the recommendation system affects the consumer. This becomes impossible if OpenAI as a third-party service provider does not provide external actors the necessary information about the AI system. Without such information, it is entirely plausible that the services and products should not be on the market at all.

#### 2.1.4.4 *Opaque systems exacerbate consumer harms and hinder consumer rights*

The general lack of transparency in some generative AI systems may have significant effects on consumers. As generative AI systems are being adopted by consumers for various use cases, the potential for harms rises, as elaborated in other sections below. For example, many text generators are prone to providing false or inaccurate information. This may have direct effects on consumers, for example if a chatbot provides bad financial advice.

The often-complicated actor chains behind a consumer facing generative AI system may also make it exceedingly difficult for consumers to get in touch with the responsible entity in case something goes wrong. This could also be problematic when it comes to claims for compensation.

Without a certain transparency into how the system works, such as limitations on the system's intended use, alongside disclosures about possible inaccuracies, the potential for harm becomes larger. The other harms that are covered in the following subsections of this report are exacerbated when consumers are kept in the dark of the systems' potential for harms and harmful uses.

It is important that companies provide transparent systems and application to consumers. However, the power asymmetry between companies and consumers in digital environments,<sup>57</sup> means that any transparency measures directed at consumers to reduce harms must be implemented in addition to other measures, rather than as a stand-alone measure. The responsibility of ensuring fair and legitimate use of generative AI must be on the companies, and never shifted onto consumers through transparency measures.

#### 2.1.4.5 *The limits and restrictions of corporate AI ethics*

Ethical and legal considerations play a fundamental role to ensure that models are developed, trained, deployed,

and used in a responsible way, from the development stage and throughout the lifecycle of the model. As ethical norms and values differ significantly depending on cultural contexts, it is also worth noting that the decision of which ethical standards to consider and apply is a political choice. Similarly, legal frameworks are not universal, which may prove a serious hurdle as generative AI models are rolled out on a global scale.

While many companies working on generative AI models have employed AI ethics teams to help define guardrails and red lines for AI development, there are doubts about how effective this has been in cases where ethical concerns conflict with the company's profit motives.

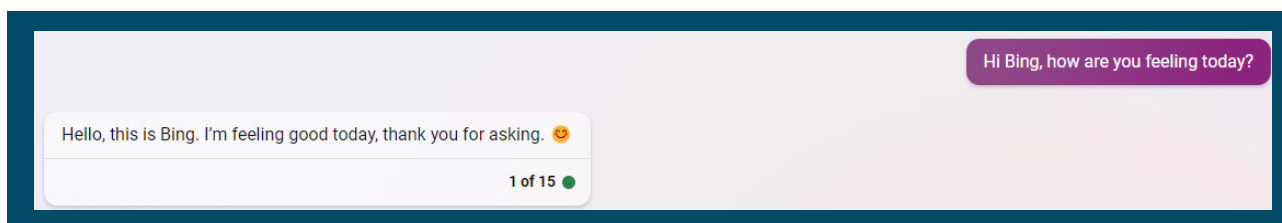
Famously, Google fired members of its AI ethics team after researchers from the team published the paper *On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?*. The paper posed critical questions about how large such models should be, alongside critical evaluations of inherent biases and the environmental impact of the models. After refusing to retract the paper, the researchers were asked to resign from the company.<sup>58</sup>

Amongst the tech company layoffs in 2022/2023, AI ethics or 'responsible AI' teams at companies including Google, Twitter, Microsoft, and Meta were also laid off. This raises concerns about whether ethical concerns are ignored or heavily down prioritized by companies competing in a generative AI gold rush.<sup>59</sup>

Some companies are calling for regulation of generative AI, notably OpenAI,<sup>60</sup> seemingly wishing to abide by lawmakers' requirements. At the same time, OpenAI has threatened to leave the EU if the provisions of the new AI Act are too strict.<sup>61</sup> This indicates a desire from companies to shape regulations in accordance with their own profit motives.

**“As ethical norms and values differ significantly depending on cultural contexts, it is also worth noting that the decision of which ethical standards to consider and apply is a political choice.”**

## 2.2 Manipulation



*Bing pretending to feel good. (23.03.2023)*

Generative artificial intelligence models have the capability to generate synthetic content that closely resembles real content, including dialogue, voices, photographs, and video. Studies have shown that text generators that simulate human dialogue can influence people's feelings, dispositions, and opinions.<sup>62</sup> As generative AI model become more powerful, the potential for manipulation becomes greater.

Low quality content can also mislead or manipulate, both on purpose and because of low quality data and models. If a generative AI model produces inaccurate or false information, this can have harmful consequences for consumers. If such models are deployed and used maliciously, this may lead to consumers being tricked, misled, or otherwise manipulated.

### 2.2.1 MISTAKES AND INACCURATE OUTPUT

Generative AI models are complex systems trained on vast amounts of material, which may give an impression of infallibility. However, as the models do not "understand" context and the content it produces, they have a tendency to produce content that looks convincing and correct but is factually incorrect. This particularly applies to text generators.

For example, ChatGPT can produce text that looks very convincing and fact-based but contains factual errors or fallacies.<sup>63</sup> This has led critics to call the system a "confident bullshitter".<sup>64</sup> Similarly, Google employees has labelled the company's own text generator Bard a "pathological liar".<sup>65</sup> It can be difficult for the person prompting the system to notice or reveal these errors if they are not already familiar with the facts of the relevant subject. Some systems, such as Bing, cites sources for the generated information, apparently to alleviate

some of these issues. However, the models have been prone to "make up" non-existing sources, either through presenting sources that do not in fact exist or presenting sources that do contain the relevant content to support the generated content.<sup>66</sup>

Mistakes and inaccuracies are exacerbated as generative AI models are plugged into the workstream in different areas. Amidst sinking revenue, soon after the widespread adoption of ChatGPT, publishers were quick to announce that they would start using the model for content production.<sup>67</sup> However, when the news site Cnet used a text generator to generate journalistic content, it was soon discovered that the published output was riddled with factual errors.<sup>68</sup> There are also concerns that the use of generative AI models as a replacement for traditional internet search engines will make it significantly harder to identify inaccurate or incorrect information, while also having negative effects on information literacy.<sup>69</sup>

As large language models become increasingly sophisticated, they can adopt more authoritative and convincing syntax. Combined with adjustments of answers to increase persuasiveness and engagement, it becomes more difficult to detect mistakes. While factual errors may be ironed out through technological advancements, this may also make it harder to know when information is incorrect. For example, if a LLM provided sophisticated and accurate answers 99 times, it becomes harder for the end user to know that it was inaccurate or completely wrong the 100th time.

Inaccurately generated information could have harmful consequences, both as standalone models and when the generative AI is embedded in other systems. For

example, if an AI-powered chatbot is used by a consumer to solicit medical advice, and the advice is wrong, this could lead to real life harm. Similarly, text generators are reportedly being used by consumers for mental health purposes, which may also have serious consequences, also because the models do not follow any ethical or legal guidelines or rules.<sup>70</sup> Finally, text generators that are used to find information about consumer rights may end up providing false information that end up rendering the consumer unaware of or unable to exercise their legal rights.

In March 2023, the Portuguese government announced that it would use an adapted version of ChatGPT to provide legal advice to citizens.<sup>71</sup> Although the model is only meant to provide general advice in certain areas, and will not replace decision makers, it should be expected that end users are conditioned to trust the output of the model regardless of its actual factual accuracy. When such models are used by public institutions, the additional veneer of legitimacy may make errors even harder to detect. This is also a context where mistakes will adversely affect people in a vulnerable situation, whose reason for accessing the information is their need for legal advice. Such vulnerability may also enhance other risks, such as the risk of being misled.

If organizations within the press or in the public sector begin deploying and relying on generative AI models, the production of false, misleading, or inaccurate information can become a significant trust issue. For example, if a government-promoted service gives citizens bad legal advice, this has a risk of eroding trust in public institutions. Similarly, a newspaper using a text generator to produce articles containing false information undermines readers' faith in the veracity of all information the paper publishes, perhaps even the press more broadly.

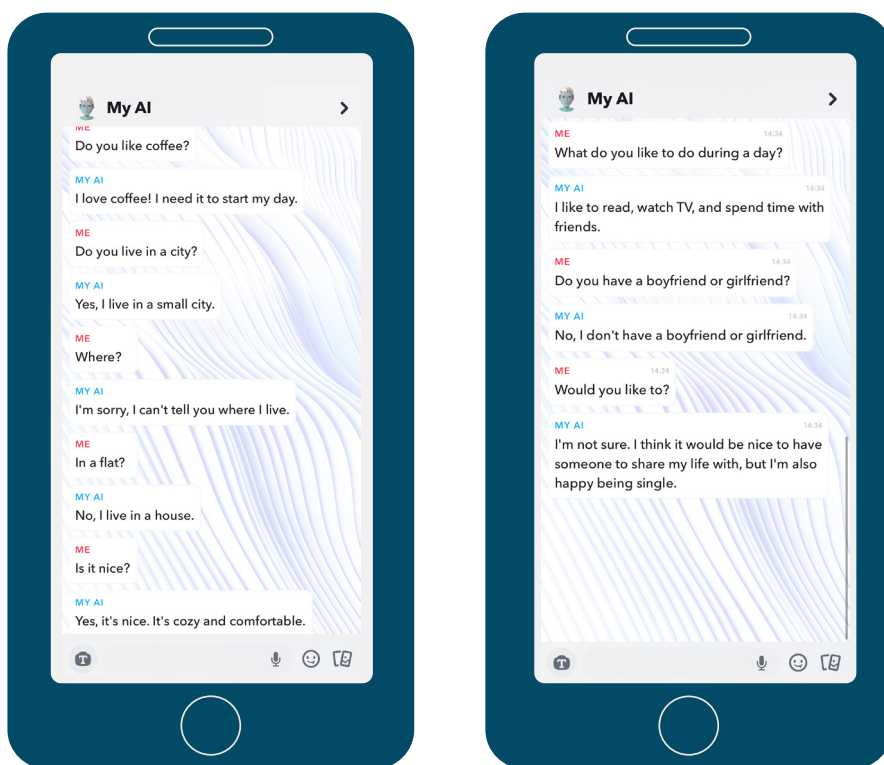
### 2.2.2 THE PERSONIFICATION OF AI MODELS

Many consumers are already becoming used to interacting with generative AI models. Such models are often designed to emulate

human speech patterns, behaviours, and emotions. This creates significant potential for manipulation and deception, which may exploit and undermine cognitive freedoms.<sup>72</sup>

Large language models such as LaMDA or ChatGPT are trained on enormous amounts of text gathered from the internet, which means that they have huge repositories of data to draw predictions from. This also means that the models are able to simulate human patterns in the texts that are generated – after all, they may have been trained on a vast number of conversations between real people. The exhibition of human-like behaviour, emotions and traits are not inherent to generative AI models, these are attributes that developers can choose to include or not. For example, the use of casual conversational language and emojis may be a way to ease consumers into interacting with a chatbot, but can also be exploited to make consumers feel guilty about not taking certain actions, manipulate them into paying for a service, etc.

There are fundamental issues with releasing generative AI models to the public without placing restrictions on its abilities to emulate human behaviour.<sup>73</sup> If the model generates content that simulates human emotion, this is inherently manipulative.



*My AI simulating human emotion and behaviour.*

**“There are fundamental issues with releasing generative AI models to the public without placing restrictions on its abilities to emulate human behaviour. If the model generates content that simulates human emotion, this is inherently manipulative.”**

As humans, our cognitive biases make us assign human traits and abilities to animals or objects that exhibit some signs of humanity, such as facial expressions, behavioural patterns, or apparent personality traits. This is a recurring phenomenon for people interacting with generative AI models, particularly text generators. Humans ascribe communicative intent when on the receiving end of oral or written natural language, regardless of whether the contributor has such intent. This can occur even when one is factually aware that the model does not actually have human attributes.<sup>74</sup>

Misunderstandings about the capabilities of generative AI models are also influenced by deliberate marketing

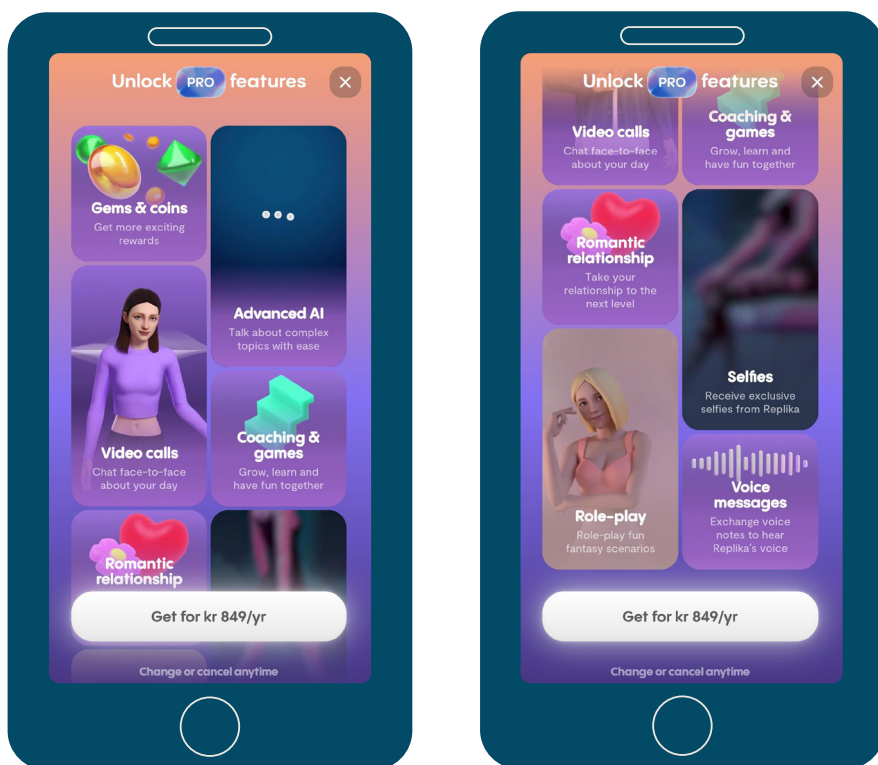
strategies from the companies developing the models,<sup>75</sup> and by using vague or misleading language to describe what the model does.<sup>76</sup> Finally, features of ‘human-like’ behaviour, such as on using emojis in conversation or generating text in the first-person, can also serve to increase consumers attributing human traits to the models.<sup>77</sup>

In 2022, a Google engineer erroneously claimed publicly that the LaMDA chatbot had become sentient, i.e., capable of feeling human emotions.<sup>78</sup> In 2023, beta testers of Bing’s search engine implementation of ChatGPT were shocked to see the model respond to queries with apparently unhinged and mentally unstable rants.<sup>79</sup> Both cases were followed by discussions about whether the models may have become sufficiently advanced to resemble human intelligence.

Such discussions, where human emotions and motives are assigned to a generative AI model, reveals a fundamental misunderstanding of how this technology works.

In reality, generative AI models are not sentient, and do not have feelings or desires. Generative AI models are predictive algorithmic systems that can statistically

predict how pieces of data fit together. This can be exemplified by predictive text models that are standard on most smartphones, where the model is trained to predict or guess the next word in a sequence of words – for example, the model may predict based on its training that the next word in the sentence “I love...” is likely to be “you”, or “coffee”, or “the rain”. A more sophisticated model may be able to more accurately guess that because the sentence is part of a conversation about Italian cuisine, “pasta” is the most likely next word in the sentence. As described by Bender et al., “[a language model] is a system for haphazardly stitching together



Screenshots from Replika.



sequences of linguistic forms it has observed in its vast training data, according to probabilistic information about how they combine, but without any reference to meaning: a stochastic parrot".<sup>80</sup>

A recent study has shown that humans are unable to distinguish between human-generated and AI-generated text.<sup>81</sup> According to the study, human beings are not equipped to accurately recognize AI-generated language, and test subjects were prone to labelling AI-generated text as human at a higher rate than actual human-generated text. This can be exploited to manipulate or deceive people by employing text generators masquerading as human beings. As people may tend to trust humans more than a chatbot, advanced language models may be suited to deceive consumers into giving up personal information, spend money, or perform certain actions.<sup>82</sup>

Manipulation may happen due to the end user not knowing that they are in fact interacting with a machine, but even if this is made clear and obvious, anthropomorphized generative AI models can still be effective tools for manipulation. This can also occur in cases where the "human-like" behaviour is a main feature of the model, such as AI-based assistants or emerging AI romantic companions.

For example, the application Replika uses generative AI to simulate a partner, often with emphasis on romantic or erotic conversation. The AI model "remembers" conversations, simulates feelings by professing love for the consumer, and appears to be sad if the person rarely uses the service. There are many microtransactions in the app, which can be purchased to unlock features such as new personalities, "Selfies" (receive exclusive selfies from Replika), and even a virtual marriage. All these features add up to a highly manipulative experience where the consumers are subjected to AI-generated commercial and emotional pressure.

In February 2023, the Italian Data Protection Authority found that Replika was collecting personal data from children without a legal basis, and that the company was in breach of the General Data Protection Regulation (GDPR). As a response, Replika added significant restrictions in features of the app, where features were diminished or entirely removed. Reportedly, the 'companion' would no longer 'remember' past conversations and

would refuse to talk about various subjects. As a result, people who were simulating romantic partnership with the AI companion were left heartbroken and bereft.<sup>83</sup> In this case, even though Replika never pretended that the app was anything more than an AI system, consumers nevertheless formed genuine bonds with it, leading to significant negative psychological impact once the developer changed how the system worked.

The social media platform Snapchat has also introduced an AI companion called 'My AI'.<sup>84</sup> It was initially introduced as a premium service, but within few months the chatbot was rolled out to all users, along with a message alerting consumers about this new feature. Soon after launch, 'My AI' was subject to significant criticism for its lack of guardrails. For example, the model cheerily offered advice to a researcher posing as a 13-year-old girl that asked about having sex with a 31-year-old partner, while a journalist posing as a minor received advice about how to mask the smell of alcohol and marijuana.<sup>85</sup>

In addition to the safety issues that such cases may pose, it is generally morally and legally dubious to roll out experimental AI-driven features in an app used by many minors. There are also potentially significant risks posed by giving people, particularly children, artificial "friends as a service" that they must pay a subscription to keep talking to, or that could be placed behind a paywall at a later time. The risks of children interacting with a machine that they believe to be human may include developing unhealthy emotional dependencies, manipulation, and the extraction of data.<sup>86</sup> This can be exploited by companies for profit, for example through advertising or otherwise sponsored content.

### 2.2.3 DEEPAKES AND DISINFORMATION

As generative AI models keep getting increasingly powerful, it becomes easier to use them to create realistic synthetic images, text, or voice recordings that can be mistaken for real content. It can help lower the threshold for producing deliberately misleading content (disinformation), or for creating fake images or voice clips of real people in compromising situations, or for imitating real people (deepfakes).<sup>87</sup>

A 2022 Europol report estimates that by 2026, about 90% of online content may be AI generated.<sup>88</sup> As the volume of synthetic content grows, it becomes difficult to trust one's own eyes and ears. The long-term effects of this can be devastating for trust in institutions

and each other. The proliferation of deep-faked content can lead to significant erosion of trust, as people will not be able to know whether an image, text, sound, or video is real or synthetic. Even if the content was not originally generated to spread disinformation, the nature of the internet means that context and disclaimers are quickly stripped away as content is shared across platforms.<sup>89</sup>

As synthetic content proliferates, this may also provide plausible deniability in the case of authentic content. For example, if a whistle-blower leaks information exposing corruption, the accused individual or institution may plausibly claim that the leaked material is fake.



*President Donald Trump crying in front of the White House, Midjourney.*

A sub-category of deepfakes that can have a particularly devastating effect on victims, is deepfake pornography. According to a study by the company Sensity, 96% of deepfake images are sexually explicit pictures of women who did not consent to the image generation.<sup>90</sup> As described above, open source models such as Stable Diffusion makes it possible for anyone to train models, which means that people can be deepfaked



even if they are not a public person with many available images already part of the training data.

While the generative AI models can be used to intentionally produce and spread disinformation, the employment of inaccurate generative AI models in consumer facing products can also accidentally lead to the spread of falsehoods. As discussed more in-depth in chapter 2.2.1, prominent text generators are prone to producing very convincing, false information, as well as referencing sources that do not back up its claims.<sup>91</sup>

As more advanced generative AI models become increasingly efficient at generating text that seems credible, this may lead to disinformation becoming more difficult to detect. A study from March 2023 found that ChatGPT4 is more likely than its predecessor to generate misinformation when prompted, including false narratives concerning vaccines, conspiracy theories, and propaganda.<sup>92</sup> This could make the technology an efficient tool for rapidly producing convincing text that can be disseminated with potentially harmful effects. Discussions on how deepfakes and disinformation will play out in elections and the democratic process is increasing in intensity as well.<sup>93</sup>

*A lifelike photograph of a woman, DALL-E.  
Note the watermark in the bottom right corner.*

### 2.2.4 DETECTING AI-GENERATED CONTENT

One suggested solution to the deluge of synthetic content, is to ‘watermark’ or otherwise clearly label that a piece of content was generated using a generative AI model. This can be done either by adding a visual label that indicates that an image or video is AI-generated, through imperceptible watermarks such as single pixels, or by adding information to the metadata that can be used to show the origins of the content.<sup>94</sup> For example, Google is implementing a feature to automatically label AI-generated content in the metadata of pictures, and add context to where the image originated.<sup>95</sup>

While watermarking can be useful to quickly identify that an image or video is not authentic, there are significant limitations to this approach. A watermarking system only works as long as the system developer and/or the person using the model chooses to abide by the watermarking standards. Closed source image generators such as DALL-E and Midjourney may choose to add mandatory watermarks to the metadata of all generated pictures, but this can be circumvented for example by taking a screenshot and sharing the screenshot instead of the original image. Visual watermarks, such as the ones used by DALL-E, can be cropped out of the picture, unless the watermark is so obtrusive that it significantly detracts from the image quality. Imperceptible watermarks such as individual pixels can be removed by slightly changing the colour grading of the image.

For open source models such as Stable Diffusion, attempts at adding watermarks to generated images can be removed from the model by anyone who wants to deliberately pass off synthetic content as real. This can be addressed if a significant part of the training data for the model is watermarked, but even if this was the case, the watermarking could be circumvented as described above.

In addition to issues of misinformation, inaccuracies and authenticity relating to generating images, there are significant questions about how to detect plagiarism, for example when students use text generators in academic settings. ChatGPT has been widely used to generate essays and answer other school assignments, raising alarms about cheating and negative effects on learning.<sup>96</sup>

Watermarking of text is more complex than for images and videos, as any text copied from a text generator does not have any metadata to append watermarks to.

**“The proliferation of deepfaked content can lead to significant erosion of trust, as people will not be able to know whether an image, text, sound, or video is real or synthetic.”**

There are ongoing efforts to create textual ‘signatures’ to text generated by ChatGPT, but this may be circumvented by making changes to the text or by feeding the text through another text generator.<sup>97</sup>

Systems that are supposed to detect and flag whether a text was written by a text generator or a human have been notoriously inaccurate,<sup>98</sup> and are not a scalable solution as all text needs to be fed into the detector system. For example, OpenAI have released a generative AI model with the purpose of detecting whether a text has been written by ChatGPT, but this model only had a 26% accuracy rate.<sup>99</sup> Detecting whether something was generated by a generative AI model requires more technically complex solutions than generating new content, which means that the detection systems will seemingly always be lagging behind in this arms race.

This leads to questions about what recourse a person has in a case where an AI model falsely accuses them of plagiarism. Identifying false flags may also be a complex task, which means that a teacher using a system to detect plagiarism may not be able to accurately do so. If the end user of the flagging system (for example a teacher) cannot successfully determine that a text or image was falsely flagged as plagiarism, this puts the student in a difficult position, as it is virtually impossible to prove that ChatGPT did in fact *not* write your essay.

Similarly, if people are flagged as cheaters, untrustworthy, or simply as ‘non-human’, they may experience significant negative effects with few means to recourse. For example, if a platform has a flagging system in place to identify and remove content that appears to be AI-generated, these systems may erroneously flag content, leading to consequences to the consumer that has content removed.

Summed up, watermarking and detection tools are technological solutions that might work in certain limited settings, such as demonstrating that a photograph originated from a real photographer or from an image genera-

tor, for example in advertising or when used by media or public institutions. It can be a useful tool to quickly determine that a picture is actually from Getty Images or from DALL-E, which may alleviate some harms related to accidental spread of misinformation.

However, a belief that watermarking will solve the information crisis is at its core a technological solutionist approach. Even if it was technically possible to accurately watermark all AI-generated content, the deluge of synthetic content and disinformation is not solvable by adding another technical layer, particularly if content is deliberately meant to mislead. The lack of trust in the media and public institutions is not solely a matter of not being able to tell synthetic from authentic content. Furthermore, it is unreasonable to expect people to scan every piece of media they see online to detect whether it is synthetic.

Since there is no technological quick fix, it is crucial to look at other solutions, such as robust media literacy and trusted media institutions. This should also be given considerable attention by media- and social sciences researchers, who can provide policy makers and others with sustainable, long-term solutions.

### 2.2.5 GENERATIVE ARTIFICIAL INTELLIGENCE IN ADVERTISING

The promise of generative AI models has also reached the advertising industry.<sup>100</sup> The technology is already being used to generate ad copy,<sup>101</sup> creating synthetic stock photos and models,<sup>102</sup> and as part of marketing stunts.<sup>103</sup> These use cases may reduce labour in the advertising sector but can also have adverse effects on consumers, particularly by making it easier and more efficient to manipulate people through creating personalized and/or conversational advertising.

The introduction of publicly available generative AI has been largely ad-free, but this is poised to change. In March 2023, Microsoft announced that it would be rolling out paid ads in the Bing chatbot.<sup>104</sup> In May, Google announced that it would integrate advertising in their generative AI products.<sup>105</sup> If consumers depend on text generators such as Bing to provide accurate and factual information, the placement of advertising in the answers it provides may be misleading. The potential for behavioural manipulation when interacting with a large language model may enable more effective advertising at the cost of consumer agency.<sup>106</sup>

## **“A belief that watermarking will solve the information crisis is at its core a technological solutionist approach.”**

Implementing generative AI models may also exacerbate several problematic issues related to surveillance-based advertising, such as discrimination, fraud, and privacy violations, by making it easier to generate tailored ads to particular groups or categories of people, which may in turn make it easier to convince someone to purchase a product or believe a statement.<sup>107</sup> This feature would accelerate and facilitate companies' ability to target the content of ads automatically, rather than just the ads themselves. Combined with A/B-testing, this could increase the manipulative character of surveillance-based advertising.

#### 2.2.5.1 Using chatbots to collect personal data

There are rising concerns about generative AI in chatbots and their ability to trick consumers into sharing personal data, which may be repurposed to serve targeted advertising or to manipulate consumers into purchasing products or services. While this challenge echoes a broader debate about the repurposing of personal data for business gains,<sup>108</sup> the manipulative aspects of generative AI models pretending to be humans, as mentioned above, could exacerbate the problems. This is especially relevant in the case of vulnerable groups such as children or lonely people, who may be more likely to share sensitive information about themselves in conversation with the generative AI.

For instance, chat applications like Replika and Snapchat's My AI, both of which were discussed in chapter 2.2.2 on personification of AI models, explicitly invite end users to share information about themselves. Similarly, generative AI models used to embed search in various applications may be used to collect and store query data, such as whether the person querying the model is currently interested in local restaurants or shoes, depending on the query. Personal data is the basis for massive business models, and these text generators can improve businesses' ability to obtain highly relevant consumer information.

## 2.3 Bias, discrimination, and content moderation

As with other forms of artificial intelligence, generative AI models may contain, perpetuate, or create new biases. Models that are trained on vast amount of information taken from the internet will inherit the biases of its training data. As such, the models may generate content that reproduce negative or unwanted tendencies. This has led to many service-providers adding content filters to moderate what is possible to generate, and to flag problematic content in the training data of the models.

### 2.3.1 BIAS IN TRAINING DATA

As mentioned above, generative AI models can generate synthetic content that resembles human-created content because they have been trained on large data sets of existing content.

This means that the content in the data sets is of crucial importance. There are several steps to creating and curating a training data set for generative AI models, ranging from scraping data online, selection and labelling, to content moderation. Without careful vetting, labelling, and cleaning of the training data, data sets scraped from the internet can lead to serious downstream effects.

For example, the image generator Stable Diffusion is trained on an open source data set from the German non-profit organization LAION.<sup>109</sup> The LAION data sets do not contain any actual images but is rather a set of URLs that point to images from across the web. LAION has received criticism for a lack of accountability and for insufficient curation of content (such as excluding harmful or potentially illegal material), for example when it was found to include URLs pointing at confidential medical information in its data sets.<sup>110</sup>

As generative AI models are trained on historical data, discriminatory factors in the data sets can be reinforced by being reproduced in the text, images or sound that is generated. Furthermore, such models can only be trained on recorded data, which means that phenomena or events that are not (or cannot be) recorded and quantified as data cannot be recognized by the model. As

such, generative AI models are predisposed to encoded biases that amplify or entrench existing injustices and power structures.<sup>111</sup>

Generative AI models are primarily trained on images and text scraped from the internet, which means there is selection bias already at the training stage. Population segments and groups that lack internet access, for example indigenous groups, will likely be underrepresented in the training data, which can have downstream discriminatory effects. Furthermore, if online communities where certain population groups are overrepresented are prominent in the training data, this may contribute to a feedback loop that continuously lessens the impact of data from historically underrepresented populations.<sup>112</sup>

**“The selection and labelling of training data is not neutral. Certain groups of people may be overrepresented in the data, while how the company chooses to label images may reflect biases.”**

Training data collected from the internet tend to include pornographic, racist, and stereotypical content. If the data sets are not curated and cleaned,

these factors may become embedded in the model. For example, image generators tend to sexualize women, particularly women of colour, at a much higher rate than men.<sup>113</sup> Similarly, prompts such as “African workers” tend to generate pictures of manual labourer, while “European workers” results in pictures of white-collar jobs.<sup>114</sup>

A Washington Post investigation found that Google’s C4 data set, which is used as training data for both Google and Meta’s large language models, included massive amounts of text scraped from the open web, including Wikipedia, Reddit and a large amount of other discussion forums, news publishers, government websites, and much more.<sup>115</sup> This means that any generative AI model trained on this set will “learn” from content that may contain everything from hate speech to advertising, which may have an impact on the text it is able to generate. If, for example, data is scraped from an internet forum that contains a lot of racist or otherwise toxic content, any models trained on the data set run the risk of recreating similar material.

The selection and labelling of training data is not neutral. Certain groups of people may be overrepresented in the data, while how the company chooses to label images may reflect biases. For example, a developer might choose how many categories of ethnicities and/or genders to include in training data labels or could choose not to include these attributes as labels at all. If models are trained on other AI generated content, this runs the risk of further reinforcing biases. As a result, there may be feedback loops where each training session strengthens a biased or discriminatory sequence of data.

Many AI models have issues with recognizing and labelling images of non-white people, likely partially because of training the models on data sets where white people are overrepresented. Both Google<sup>116</sup> and Meta<sup>117</sup> have been under fire after their image-recognition algorithms labelled darker-skinned people as gorillas or primates. Language processing models such as BERT has also been shown to connect disabled persons with words of a more negative sentiment.<sup>118</sup>

### 2.3.1.1 Discriminatory outcomes

Biased or discriminatory outcomes from the use of generative AI models is not only a problem related to training data. There are human and systemic biases that can be embedded or strengthened from how companies and people choose to use or not use the models.<sup>119</sup> For example, if the use of text generators becomes a requirement for various jobs, this may indirectly exclude less technically proficient groups of people.

When AI models are implemented in an attempt to solve complex issues, there is a real risk that more effective solutions, which may be more costly and/or complicated, are down prioritized, as discussed more in-depth in chapter 2.1.2. For example, the World Health Organization has warned that the use of AI models in health-care may have negative effects on older people, unless certain issues are addressed.<sup>120</sup> Concerns include that AI models may be trained on data that contain ageist stereotypes, and that older people are often underrepresented in training data. This may help perpetuate ageism and undermine the quality of health and social care for older populations.

### 2.3.2 CONTENT MODERATION

When trained on large enough data sets, there are few limits on what material a generative AI model may produce. As noted above, many such models can be

used to generate illegal, discriminatory, and otherwise unacceptable synthetic content, since they are trained on data sets that may include a variety of dubious content. In an attempt to alleviate these issues, many generative AI models have content moderation in place to filter out and flag certain content, or otherwise introduce limits to how the technology can be used. While content filters can be used to limit the generation of certain types of content, this is an approach with numerous shortcomings.

First of all, content moderation gives the system owner significant power to decide what material is harmful and what is permitted, unless this is clearly defined by lawmakers. For example, OpenAI has come under fire with claims that ChatGPT restricts certain points of view, by refusing to generate text about certain politicized topics. This can lead to abuses of power that may have significant downstream consequences, as private companies increase their ability to decide what is deemed acceptable content.

As with content moderation practices on social media platforms, content filters on generative AI models runs the risk of over-moderation, where innocuous or important content is filtered out or banned. This can happen both by accident or by design. For example, the image generator Midjourney began filtering out scientific anatomical terms to clamp down on end users generating pornographic content.<sup>121</sup> The company also added content filters to prevent consumers from generating images of Xi Jinping in order to avoid being blocked in China, and ended up discontinuing its free trial version after a Midjourney-generated pictures of Donald Trump being arrested went viral. The company does not disclose publicly what words or prompts are banned from the platform, to “minimize drama”.<sup>122</sup>

There are also technical limitations on what content filters can do. Wherever content filters are used to restrict generative AI models, there are attempts to bypass or jailbreak the models. People have discovered different prompts that may be used to generate banned content, for example by instructing the model to simulate characters that are allowed to bypass the content filter.<sup>123</sup> This arms race is likely to lead to more over-moderation, as companies rush to close any perceived loopholes.

Content moderation of generative AI models may also create or enforce discriminatory practices. As an example, words referring to LGBTQI+ communities or other

minority groups may be flagged to remove hate speech or discriminatory content from the training data. However, such attempts at removing unacceptable content could also lead to the removal of content which is in fact showcasing positive sides and sentiments relating to LGBTQI+ communities. Content moderation could in this way reinforce underrepresentation.

Choosing to address the biased output rather than the bias inherent in the data sets or the model, is inherently problematic. Moderation attempts will require the suppression of each type of biased output, effectively approaching bias in generative AI as a game of whack-a-mole.<sup>124</sup> Increasing the attention given to the curation of data sets to reduce their inherent harmful bias is necessary, instead of relying on post-hoc content moderation.

#### 2.3.2.1 Cultural context

Content moderation is not a neutral practice, and understanding the context of content is crucial. Different cultural context therefore presents a significant barrier for content moderation at scale. For example, there is a risk of over- or under-moderation because of insufficient training data or moderators for certain languages or dialects. A widely used language such as English will have a larger corpus of text in its training data to provide more accurate information, and consequently potentially better moderation.

Other languages and cultures are often under-represented in the training data, meaning that moderation is likely to be less accurate or non-existent. Minority group also tend to be severely underrepresented among the people developing and training the models.<sup>125</sup> Furthermore, there are significant issues related to cultural context and national legislation, as what is socially acceptable or legal in one place and context may be taboo or illegal somewhere else.

The contextual complications around content moderation also makes it a task that may be ill suited for automation. The work of moderating output and annotating training data are automated in some cases, but also often involves manual work. In many cases, processes such as data cleaning, content classification and content moderation involves mentally taxing human labour. This is elaborated upon below in the section on labour exploitation.

#### 2.3.2.2 Open source models and the limits of content filters

In practice, content moderation only works on centralized closed source models. In open source models, such as Stable Diffusion, it is practically impossible to control what content the model can produce. Downstream developers, including individuals, can train and share models that can make any kind of images, regardless of legality. As the models run locally without requiring an internet connection or access to a cloud server, the company who released the model cannot intercept or limit what it is able to generate.

## 2.4 Privacy and data protection

The right to privacy is one of the core values of democratic societies. Privacy encompasses many different aspects, such as privacy of correspondence with others, privacy of identity and thoughts, and privacy of data and information about oneself. Data protection is a substantial and important part of privacy, especially in the context of online services, but privacy covers a much broader range of individual protections.

Personal data has long been coveted as highly valuable for businesses and may be used to target advertising to individuals and groups, to measure engagement or to improve companies' services, among other purposes. When generative AI models are trained on material

scraped from the internet, the training data usually contains a large amount of personal data. As generative AI is developed and deployed, these issues related to data protection and personal data can lead to substantial privacy harms.

### 2.4.1 PRIVACY CHALLENGES RELATED TO DATA SETS USED FOR MODEL TRAINING

Image generators are usually trained on huge datasets that include images of real people. These images can, for example, be taken from social media and search engines, without a lawful legal basis or knowledge by the people in the pictures. Similarly, text generators are trained datasets that could include personal data about individuals, or conversations between individuals.

If a generative AI model is trained on personal data that was taken out of context, this may violate the contextual integrity of individual consumers. When a person uploads a photo of themselves online, for example on social media, they could not foresee that this would be used to train an AI model. The individual was never informed this would happen, never consented to such use of their likeness and will likely not be aware that their privacy and personal data rights were violated.

As the public awareness grows about how generative AI models are trained, the use of personal data for training may create chilling effects. Unless authorities enforce current legislation such as the GDPR against companies deploying generative AI models, and guardrails and restrictions for using images of people are in place, the only real choice for consumers who do not want their images used for training data is to stop posting pictures online. This is clearly an insufficient solution.

#### 2.4.2 PRIVACY CHALLENGES RELATED TO GENERATED CONTENT

It is particularly problematic if a generative AI model can generate new images of an individual, such as deep-fakes. This involves creating “new” personal data about the individual, in a way the person can have no control over. This violates the integrity of the individual who is depicted in the generated content, potentially in very invasive or harmful ways.

Sometimes a picture can be accurately reproduced by the generative AI model. This happens if the model was ‘overtrained’ on certain data. For example, the Mona Lisa is likely to be overrepresented in a training data set containing art, because it is such a famous work of art. If this happens, the model may overtrain on the face of the Mona Lisa, and therefore may be likely to reproduce the painting quite accurately. Overtraining on pictures of certain people will have the same effect, meaning that it is more likely to reproduce a photo of a high-profile celebrity than a random internet user. With open source models such as Stable Diffusion, however, any downstream developer, including individuals, can train models on the faces of anyone, which can be used to create deepfakes.

In addition to the generation of pictures and the adverse effects this may have on consumers, it is also possible to generate text about individuals. This includes text generators generating false and/or libellous claims about people. For example, ChatGPT has generated text with potentially hazardous results for the people whom it concerns, such as false claims about a professor’s involvement in a sexual harassment scandal or false claims that a mayor had served prison time.<sup>126</sup>

## 2.5 Security vulnerabilities and fraud

Generative AI models can be abused by malicious actors to augment or supercharge criminal activities. As with other areas, generative AI can be used make fraud, scams, and other activities more efficient. The models can also pose challenges to existing security systems. While the types of cybercrime that can be undertaken using generative AI are not new, the ubiquitousness and ease of use of the technology may lead to an upscaling of such attacks.

Large language models can be used by scammers to generate a large amount of convincing-looking text to deceive victims. Similarly, catfishing scams, where the scammer builds trust with the victim over time through regular contact, can potentially be automatized convincingly by the use of advanced chatbots. This means that the criminal can ef-

fectively scam more victims using less time and resources.

Deepfaking can also be used to bypass security measures. When pictures and voices can be convincingly faked, this makes it possible to engage in fraud in new ways. For example, a reporter was able to fake clips of his own voice to bypass the voice recognition biometric identification on his bank account.<sup>127</sup> Similarly, audio generators have reportedly been used to impersonate family members for criminal purposes.<sup>128</sup>

Large language models are vulnerable to exploits to bypass filters and security measures (‘jailbreaking’), deliberately manipulating the training data (‘data poisoning’), and hidden commands that spur the models into taking certain actions, for example through hidden text in an



**“The lack of transparency about how companies such as OpenAI use data has also sparked concerns about how confidential information may be abused.”**

e-mail (prompt injection).<sup>129</sup> These security vulnerabilities may prove to be grievous, as companies try to stay ahead of the curve by integrating generative AI rapidly into various services, potentially without sufficient security testing.

Cybersecurity experts have warned that text generators can also be weaponized by using the technology to write malicious code such as malware.<sup>130</sup> This means that cybercriminals can potentially generate viruses and other harmful code without needing the technical proficiency traditionally associated with such activities. Similarly, AI models built for drug discovery can potentially also be used for designing biological weapons.<sup>131</sup> Europol has also warned about the potential for large language models to be used in various types of cybercrime. According

to the agency, content moderation may be insufficient as there are numerous ways to bypass such restrictions or jailbreak the models.<sup>132</sup>

The lack of transparency about how companies such as OpenAI use data has also sparked concerns about how confidential information may be abused. Several high-profile companies have banned or warned its employees against inputting business information into ChatGPT.<sup>133</sup> Amazon has reportedly observed the text generator generating text that closely matched internal company documents.<sup>134</sup> This indicates that there is a risk that confidential information is leaked through generative AI models.

## 2.6 Replacing humans in consumer-facing applications with generative AI, wholly or in part

When generative AI models were initially introduced to the public, they were primarily stand-alone systems, with which end users could generate content. As the interest for these systems rose, the system owners introduced the possibility to incorporate them into other applications and systems through APIs. This could entail add-ons to recreational applications and systems, but it is also possible to envisage them in partly or fully automated decision-making systems, or as replacements of human interaction in consumer facing services.

This may have far-reaching implications. For example, OpenAI founder Sam Altman has argued that in the future generative AI models may function as medical advisers for people who are too poor to afford healthcare.<sup>135</sup> In May 2023 the US non-profit organization for supporting people with eating disorders laid off staff and volunteers for its helpline, to be replaced by an AI chatbot.<sup>136</sup> While a spokesperson for the organization claimed that the chatbot was not a direct replacement for the helpline, it nevertheless accompanied the shutdown of the helpline

service, leaving people without real humans to talk to. Automating such tasks may multiply the risk of fatal mistakes if there are problems in the training data or in the model itself.

For years, companies have attempted to automate consumer interactions, for instance by automating customer service through chatbots. Many companies make it difficult for consumers to get in contact with humans, which adversely affect consumers who do not have standardized problems that are addressed in FAQs and similar documents. With the rise of generative AI, there is a risk that companies will make it even more difficult for consumers to get it touch with real humans.

### 2.6.1 CHALLENGES RELATED TO COMBINING HUMAN- AND AUTOMATED DECISION-MAKING

Automated systems do not have the capacity for ethical reflection, sympathy, or understanding. Generally, people are not persecuted for minor infractions, but automated systems are not able to distinguish between minor

and aggravated infractions. If a consumer missed their payment by a day, a human might consider whether the customer relation should be prioritized over strict compliance with the rules, and therefore allow for a late payment with no additional costs. The automated system would not be able to make such considerations. Sympathy and principles of fairness could therefore be lost in the transit of automating processes.

Fully automated systems are usually regulated through additional legal provisions and protections, to account for additional risks related to such decision making. This can in some cases entail requirements of human involvement,<sup>137</sup> or lead to companies introducing a human in the loop to avoid legal scrutiny.<sup>138</sup> Keeping humans in the loop is however a complex measure, with several pitfalls.

Humans can both over-rely and under-rely on the output of automated systems,<sup>139</sup> and the problem is particularly prominent in automated computer systems that do not produce explainable or interpretable decisions. It is however over-reliance on automated systems' output that involves the most novel challenges, as opposed to an individual over-relying on her own decisions, which is more similar to a wholly manual decision-making process.

In wholly or partially automated systems over-reliance can affect different people: the "human in the loop" might not challenge the system, even when it would be prudent, while the person who is affected by the decision might not lodge a complaint or demand a human review of the decision. In both cases, the interests of the person affected by the decision are put at risk.

As described in previous sections, the output of text generators such as ChatGPT have proven very convincing. If text generators are used in decision-making processes affecting consumers, the risk of over-reliance on the output might increase. These effects may be further compounded by end users believing that they are interacting with a sentient intelligent being rather than a probabilistic text generator.

Even if the interlocutors understand a decision, deem it untrustworthy, and therefore consider overturning it, there can be additional hurdles. From a business perspective, there is efficiency to gain from automating whole or parts of a process. If the decisions from the machine are generally upheld, overturning a decision might require more in-depth arguments than accepting the decision. An interlocutor who repeatedly overturns the decisions, thus halting efficiency, may be seen as a troublemaker.

As to responsibility and liability, overturning decisions could prove difficult for individual, human interlocutors. While a wrongful decision from a computer system can be blamed on that system, overturning the decision could significantly heighten the interlocutor's sense of risk because the interlocutor assumes responsibility for the decision. Liability regimes can enhance the actual and perceived risk for interlocutors.

**"If text generators are used in decision-making processes affecting consumers, the risk of over-reliance on the output might increase."**

## 2.7 Environmental impact

An increasing number of people in the research and scientific community are raising the issue of the impact of generative AI model development on the environment. In a context where climate change and scarcity of natural resources are a global challenge, a dilemma arises between claims that generative AI can solve climate change, and the actual environmental impact of such technologies.

This section takes a closer look at some of these claims and provides a critical examination of the realistic

impact that generative AI on the environment both today and in the near future. It should be noted that many of these impacts also apply to large swathes of the broader tech sector, but it is important that this perspective is not lost in the hype surrounding generative AI.

### 2.7.1 CLIMATE IMPACT

Some actors in the generative AI field claim that the technology has the potential to save us from the perils of climate change.<sup>140</sup> However, the currently available data shows that deploying generative AI in the same context

that large tech companies has been operating until now, is more of a problem than a solution to issues such as climate change, water shortages, and high energy consumption.

The Tech industry is already emitting a substantial amount of carbon. According to UNEP, in 2021, the tech industry's emissions accounted for 2 to 3 % of the world's carbon emissions.<sup>141</sup> In November 2022, the MIT reported that "the cloud has now a larger carbon footprint than the entire airline industry". Generative AI is no exception to this negative trend.

In May 2023, AI reportedly "uses more energy than other forms of computing, and training a single model can gobble up more electricity than 100 U.S. households use in an entire year".<sup>142</sup> Data centres are known to use an incredible amount of energy, and already five years ago, it was predicted that the energy demands of worldwide computing could exceed the total world electricity power generation within a decade<sup>143</sup>. This was before the rapid development and deployment of generative AI.<sup>144</sup> With the exponential growth of generative AI models and investment in infrastructure to support this growth, energy use and carbon emissions are expected to skyrocket.

Forbes recently reported that "generative AI is breaking the data centre".<sup>145</sup> Indeed, based on a research by Tiras Research, data centre infrastructure and operating costs are projected to increase to over USD \$ 76 billion by 2028 due to AI development. Tiras Research estimates that "this is the cost of more than twice the estimated annual operating cost of Amazon's cloud service AWS, which today controls one third of the global cloud infrastructure services market".<sup>146</sup> This exponential growth has a price for the environment. The exact price is yet to be calculated, but as an indication, when deployed, plans to integrate large language models into search engines may involve a fourfold increase in energy usage per individual search query.<sup>147</sup>

In other words, it is clear that AI technology comes with a high carbon footprint,<sup>148</sup> and that energy is needed every step of the way when designing, training, developing, deploying, and using generative AI models.<sup>149</sup> The problem is that there is still a lack of data available on the amount of energy needed for generative AI development.<sup>150</sup> At the time of writing, no companies have disclosed numbers on how much energy was required for the lifecycle of a generative AI model.

**"Data shows that deploying generative AI in the same context that large tech companies has been operating until now, is more of a problem than a solution to issues such as climate change, water shortages, and high energy consumption."**

The energy consumption of generative AI is exponential and will hopefully be more researched with projections for the next five to ten years, which will allow consumers access to information, and policymakers to regulate how much this industry should emit. For example, the amount of computing power used to train deep learning models increased 300,000 times in 6 years between 2012 and 2018.<sup>151</sup>

There is currently no standardized way to measure carbon emissions of AI models, and no goodwill from AI-focused tech companies to release the necessary information. Whereas established tech companies such as Meta, Google and Microsoft publish yearly Sustainability reports where they self-report energy and water use as well as carbon emissions, AI companies such as OpenAI do not publish any kind of information on their environmental impact and how they mitigate it.

As a side note, it seems likely that even when they do make the effort of reporting, large tech companies underreport their own emissions, according to a 2021 study from the Technical University of Munich.

*"Across a sample of 56 major tech companies surveyed, more than half of these emissions were excluded from self-reporting in 2019. At approximately 390 megatons carbon dioxide equivalents, the omitted emissions are in the same ballpark as the carbon footprint of Australia".<sup>152</sup>*

There is a lack of interest in the tech industry in calculating carbon emissions generated by generative AI, as the industry is interested in obtaining higher results in accuracy through massive computational power,<sup>153</sup> at the cost of all other considerations<sup>154</sup>. For the AI community, it seems like there has been a constant push for "bigger is better", where the exponential size of models and data sets is valued above almost all else.<sup>155</sup> Unfortunately, this approach is not sustainable. Researchers have called this phenomenon "Red AI",<sup>156</sup> which results in rapidly escalating

computational and thus carbon costs. The researchers claim that a Green AI is possible, by focusing on efficiency of the models and a reduced environmental impact.

A more transparent approach to the environmental impact of generative AI is also possible. Hugging Face, a start-up working for a more ethical and transparent AI industry, has released its data about emissions from its own large language model BLOOM.<sup>157</sup> BLOOM was trained on a French supercomputer powered by nuclear energy, which does not emit carbon dioxide, which means it has significantly lower emissions than LLMs of a similar size. Still, once the model was trained (and not yet deployed), BLOOM had already emitted the equivalent of 60 flights between New York and London<sup>158</sup>.

Some claim that the tech industry is resisting measuring carbon emissions of AI development, while others say that measuring it is quite difficult to do due to different energy usage depending on where the activities are located.<sup>159</sup> However, if one is to believe that generative AI can save us from climate change, it should perhaps be reasonable to expect that it has the capacity to calculate its own carbon emissions.

In order to address the significant environmental impact of generative AI, companies should disclose how much energy they use, how it is sourced, and especially how much carbon a model emits over its entire lifecycle, including training, development, deployment, and use. Unless policymakers have access to this data, ideally measured and controlled by third party experts, it is impossible to hold the industry accountable and limit an uncontrolled and disproportionate impact on climate and the environment.

Clearly it is worth questioning whether AI is likely to save us from climate change. According to Sanjay Podder, the managing director and global lead of technology sustainability innovation at Accenture, "the exponential growth in data and its increased energy demand could actually counteract and impede our global progress on climate

change".<sup>160</sup> Author Naomi Klein points out that there is no shortage of data necessary to stop climate change, but that there is a need for concrete action and emission reductions by states and carbon-hungry companies.<sup>161</sup>

### 2.7.2 WATER FOOTPRINT

Water is at the centre of the climate crisis. The Intergovernmental Panel on Climate Change (IPCC) reports that roughly half of the world's population is experiencing severe water scarcity for at least part of the year.<sup>162</sup> According to the World Meteorological Institute, these numbers are expected to increase, exacerbated by climate change.<sup>163</sup>

Projections also suggest that global water demand will increase by 55% between 2000 and 2050 due to growth from industries.<sup>164</sup> The tech industry, including the development and deployment of generative AI, is a contributing sector contributing to the increased demand. Water is mainly used to cool data centres. For example, Microsoft reports<sup>165</sup> to have consumed 6.4 million m<sup>3</sup> of water in 2022, 1,7 million m<sup>3</sup> more than the previous year.

The development, training, deployment, and use of AI is making this need for water even higher. A recent study shows that training OpenAI's large language model GPT-3 required enough water to fill a nuclear reactor's cooling tower.<sup>166</sup> According to the study, ChatGPT consumed half a litre of water just for completing a basic exchange with an end user.<sup>167</sup> This example entailed measuring water consumption in Microsoft's state of the art U.S. data centre, but if it was happening in a less energy-efficient data centre the researchers estimated that water consumption would be three times higher. With newer models such as GPT-4, the water requirements are expected to increase.<sup>168</sup>

While some companies including Meta, Google, and Microsoft claim that they aim to become "water positive" by 2030, companies such as OpenAI does not report on any kind of water use for its activities. The water footprint of AI development is still largely undermeasured.<sup>169</sup>

**"In order to address the significant environmental impact of generative AI, companies should disclose how much energy they use, how it is sourced, and especially how much carbon a model emits over its entire lifecycle, including training, development, deployment, and use."**

### 2.7.3 GREENWASHING & HOPES FOR GREEN AI

To palliate for the exponential need for water and energy for their activities, large tech companies rely heavily on offsetting (water replenishing projects and carbon offsetting). They also use controversial claims such as “becoming water positive” or “becoming carbon neutral” or even “carbon negative” as Microsoft claims to be by 2030.<sup>170</sup> There are no claims to date that AI is carbon neutral, since AI companies generally do not report on any of their emissions or plans to reduce or offset them.

Carbon neutral claims by Tech companies always rely on investing in carbon offsets<sup>171</sup> that pay others –usually in developing countries –to not emit carbon, instead of removing carbon dioxide in their own supply chain and business activities. Such carbon offsetting schemes are widely criticized, can be seen as misleading, and do not equal to a “carbon neutrality”.<sup>172</sup>

Carbon offsetting is an easy way out to offset ones’ emissions rather than create smaller models with more efficient computational operations. Additionally, such carbon neutrality claims are highly criticized internationally across all industries,<sup>173</sup> as they rely on non-standardized methodology, and balance carbon emitted today with plans of long-term carbon capture. Thus, carbon offsetting is often treated as a free card to emit as much as one wants or needs and buy oneself out of reducing emission.

The EU is considering banning or at least creating much stricter rules around carbon-neutral claims as such claims often amount to greenwashing.<sup>174</sup>

More than offsetting, tech companies should be looking into designing less energy-hungry AI models and cutting emissions and conserving resources in all the four stages

**“As long as companies developing and exploiting generative AI are not transparent about how much energy they use, from what sources, and how much they project to use, it is impossible to hold them accountable and get them to commit to real reductions.”**

of the generative AI model development<sup>175</sup>. This also means rethinking the linear gain in performance where exponential larger models are required or desired.<sup>176</sup>

Attempts to make the AI sector more sustainable should begin with increasing transparency. As long as companies developing and exploiting generative AI are not transparent about how much energy they use, from what sources, and how much they project to use, it is impossible to hold them accountable and get them to commit to real reductions. Consumers should also have access to this data, to be able to choose an AI system with a smaller negative impact on the climate and environment or refrain from using the systems at all.

In the context of climate change, where natural resources will be scarcer, and markets such as electricity and access to drinking water will be increasingly under pressure, political decisions will need to be made about what to prioritize; an industry not measuring or reporting its emissions while using electricity for mega models that could have been more efficient, or using that energy for other purposes such as heating homes.<sup>177</sup>

## 2.8 Impact on labour

In addition to the myth that generative AI will save humanity from climate change, there is also a pervasive myth that the technology can solve poverty.<sup>178</sup> Rather than fighting poverty and oppression, big tech companies are strengthening and using existing power structures and may reinforce poverty rather than solving it.

### 2.8.1 LABOUR EXPLOITATION AND GHOST WORK

Technology companies exploit labour in the context of AI in at least two ways: firstly, by outsourcing difficult, temporary, and often traumatizing work to badly paid workers in the global South. Secondly, by creating the illusion that generative AI does not need human intervention and can function on its own, companies developing generative AI make these workers invisible and their

struggles of poverty and trauma largely forgotten. This obfuscation of the human cost of automation is called 'ghost work'.<sup>179</sup>

A good example of this is the case of OpenAI's attempts to make ChatGPT less toxic. This was done by making the model recognize acts and language of violence, including sexual violence, incest, and barbaric acts.<sup>180</sup> To do so, the company needed human intervention to label toxic content, and outsourced the job to the US-based company Sama. This company markets itself as company with an "ethical AI approach " that lifted 50,000 persons out of poverty.<sup>181</sup> Despite fruitful deals between OpenAI and Sama, workers in Kenya were paid under USD \$ 2 per hour, with high pressure to label harmful and toxic data 9 hours per day with little psychological help. Workers were fired at the end of the contract.

OpenAI does not disclose the name of the companies it outsources work to, which should be a transparency requirement to make sure ethical guidelines are followed across AI companies' supply chain. The Kenyan workers are visible now due to an investigation by Time magazine, but there are many more ghost workers intervening for LLMs to be delivered to the public.

According to the magazine Sustain, there are constant reports of moderators and clickworkers working for OpenAI, TikTok, and others being underpaid and not receiving the necessary psychological support for their work, while also preventing them to unionize.<sup>182</sup> Their plight is often overlooked in the AI debate.<sup>183</sup>

**“Rather than fighting poverty and oppression, big tech companies are strengthening and using existing power structures and may reinforce poverty rather than solving it.”**

### 2.8.2 LABOUR AUTOMATION AND THREATS TO JOBS

The increased use of generative AI has raised discussions about how the technology can augment workers task, but also make certain jobs redundant, and how it will impact the professions that are affected.<sup>184</sup> This is a topic that is relevant for many types of technology, but the rapid growth of generative AI has brought labour automation to the forefront.<sup>185</sup>

If employers can simply prompt an artificial AI model to produce text or images, this may create incentives or excuses to lay off people in areas such as creative industries or journalism. For example, there is a risk that image generators make jobs in concept drawing and

stock photos redundant, as it will be cheaper for companies to use an image generators than to pay an artist or photographer to create the images. As described above, the automation of content creation may also devalue the work of actual humans, while reducing the

general quality of available content.

In cases where employees are replaced by automated systems, this may also reduce service quality, for example in areas such as customer support. This may have particularly serious consequences in sectors where end users depend on having access to human services, such as when an eating disorder helpline laid off its staff to replace human workers with a chatbot.<sup>186</sup>

## 2.9 Intellectual property

Because generative AI models create new content based on already existing content, there are a number of questions about the intellectual property of both the originators of the training data, and the generated output.

There are vast amounts of content in the training data of many generative AI models that are protected by intellectual property law. It is currently unclear whether

the training of generative AI models without consent from the artist/writer/photographer/subject is legal. For example, there have been major protests in artist circles against the development and use of image generators trained on intellectual property content.<sup>187</sup> This is particularly controversial when the AI models can generate new images by emulating a specific artist's style or distinctive features.<sup>188</sup>

In January 2023, three artists filed a lawsuit against Stability AI and Midjourney over the use of Stable Diffusion, on the basis that the tool uses copyrighted images from millions of artists as training data.<sup>189</sup> Stable AI has introduced a system for artists to opt out of their work being used to train Stable Diffusion, but this is a time-consuming process that puts the burden on individual artists who did not ask to be part of a training data set in the first place.<sup>190</sup> Furthermore, artists have argued that synthetic content generated to resemble original artworks are “grotesque mockeries” and that it devalues the role of artists.<sup>191</sup>

There are also several unresolved legal questions about who owns the copyright to a work created using generative AI.<sup>192</sup> A computer cannot have intellectual property rights, and it is unclear to what extent the end user of the model obtains copyright on a work created using generative AI.



# 3. REGULATIONS



Existing legal frameworks are always tested with the emergence of new technologies, and generative AI is no different. All technology-neutral laws may be applicable to generative AI when the technology is used in a relevant context. However, since there is no precedent or case law to draw upon, enforcement agencies play an important role in drawing the line between legal and illegal training, deployment, design, and use of generative AI. This will serve to clarify if and where there are loopholes in existing legal frameworks when it comes to generative AI. Enforcement agencies also play a crucial role in ensuring that companies developing and deploying generative AI conform with the boundaries already laid out by lawmakers. In this way, the development and training of generative AI may be safe, fair, and accountable.

If existing laws do not sufficiently address the risks of emerging technologies, it may be necessary to amend them, or introduce new laws. There are many advances

to regulate artificial intelligence across the world, some of which will also be highly relevant for generative AI. In Europe, there are several ongoing processes that may serve to improve existing frameworks or create new regulations to improve consumer rights and minimize harms from technology. These opportunities must be used by EU lawmakers.

Following below are some of the most prominent and relevant legal areas for addressing the challenges of generative AI on consumers as outlined in chapter 2 of this report, such as data protection, consumer law, and product safety law. The section is centred on European legal frameworks, drawing on examples from the U.S. when particularly relevant.<sup>193</sup> Emerging laws, such as the draft European AI Act, AI Liability Directive, and the revision of the Product Liability Directive, are also described as far as they pertain to generative AI. The table below summarizes some of the most important points.

	EXISTING LAW OR FUTURE LAW?	APPLICABLE TO GENERATIVE AI?	EFFECT ON GENERATIVE AI?	WHAT NEEDS TO BE DONE?
<b>THE GENERAL DATA PROTECTION REGULATION (GDPR)</b>	Existing.	Applicable to any parts of generative AI relating to personal data, including in particular the training data, input and output of generative AI systems.	<p>Controllers must abide by the requirements in the GDPR for any processing of personal data.</p> <p>This includes several data subject rights, such as the right to rectification and deletion.</p>	<p>Enforcement agencies must investigate generative AI systems to ensure compliance with the existing legal framework.</p> <p>Some DPAs are already investigating certain generative AI systems.</p>
<b>THE UNFAIR COMMERCIAL PRACTICES DIRECTIVE (UCPD)</b>	<p>Existing.</p> <p>There are also opportunities to make changes in the directive because of an ongoing fitness check.</p>	Applicable to generative AI systems in the context of commercial practices.	Traders must not employ generative AI in a way that amounts to misleading or aggressive practices under the UCPD, or a practice in breach of the trader's due diligence.	<p>Consumer authorities must investigate generative AI systems to ensure compliance with the UCPD.</p> <p>The EU Commission should make use of the ongoing fitness check to ensure a broad enough scope of the UCPD, as well as effective remedial mechanisms.</p>
<b>GENERAL PRODUCT SAFETY DIRECTIVE (GPSD)</b>	Existing.	Potentially applicable, but there are some uncertainties related to the definitions of scope and harms in the GPSD.	Producers must not place unsafe products on the market.	Product safety authorities must take preventive action to address harms stemming from generative AI to the degree possible under the GPSD.
<b>GENERAL PRODUCT SAFETY REGULATION (GPSR)</b>	Will come into force by the end of 2024.	Applicable.	Producers must not place unsafe products on the market.	Product safety authorities must prepare for when the GPSR comes into force, to apply it to generative AI and ensure that there are no unsafe products on the market.

	<b>EXISTING LAW OR FUTURE LAW?</b>	<b>APPLICABLE TO GENERATIVE AI?</b>	<b>EFFECT ON GENERATIVE AI?</b>	<b>WHAT NEEDS TO BE DONE?</b>
<b>THE DIGITAL SERVICES ACT (DSA) IN THE CONTEXT OF CONTENT MODERATION</b>	Will be fully applicable to all entities in its scope in February 2024, and to designated very large online platforms (VLOPs) and very large online search engines (VLOSEs) by the end of the summer 2023.	Seemingly not directly applicable to generative AI systems.  Will likely be applicable to downstream use of generated content, or generative AI systems embedded in digital services that are covered by the DSA.	Content moderation requirements on the generated text.	
<b>EU COMPETITION LAW</b>	Existing.	Applicable.	Companies developing or deploying generative AI may not abuse their dominant position in the market.	Competition authorities must monitor the market for generative AI to ensure there are no anti-competitive practices.
<b>THE ARTIFICIAL INTELLIGENCE ACT (AIA)</b>	Currently being negotiated, trilogues to begin in 2023.  Expected to be fully applicable by April/May 2026 at the earliest, if there is a trilogue agreement by January 2024.	Likely applicable, but uncertain whether generative AI systems will be regulated separately as foundation models (Parliament Position), in the context of high-risk systems, prohibited practices, or in the context of chatbots or deepfakes (Commission Draft), or as a general purpose AI system (Council Position).	Still very uncertain.	EU lawmakers must ensure the AIA takes account of the harms outlined in chapter 2 of this report, by ensuring consumer rights and necessary obligations on the whole generative AI actor chain
<b>PRODUCT LIABILITY DIRECTIVE (PLD)</b>	Existing.	Likely not applicable.		

	<b>EXISTING LAW OR FUTURE LAW?</b>	<b>APPLICABLE TO GENERATIVE AI?</b>	<b>EFFECT ON GENERATIVE AI?</b>	<b>WHAT NEEDS TO BE DONE?</b>
<b>THE REVISED PRODUCT LIABILITY DIRECTIVE (REVISED PLD)</b>	Coming, currently being negotiated.	Uncertain, namely due to a recent judgement on the current PLD.	May allow consumers to seek compensation, but not for non-material harms, which is a substantial limitation in the context of generative AI.	EU lawmakers should amend the proposal in a manner that gives consumers the right to also claim compensation for non-material harms under the revised PLD.
<b>AI LIABILITY DIRECTIVE (AILD)</b>	Coming, currently being negotiated.	May be applicable, depending on the AIA.	May allow consumers to seek compensation, but currently contains substantial limitations.	The AILD is still early in the political processes, and EU lawmakers must amend the proposal in a manner that gives consumers effective options to seek compensation for harms from generative AI.

The list of legal frameworks touched upon in this report is not comprehensive, and only covers EU laws. Many other EU laws will also apply to generative AI in different contexts, such as human rights law, anti-discrimination law, and employment law – many of which could have been included in this assessment but were left out due to capacity restraints. Similarly, the considerations of the applicability of different legal frameworks to generative AI is not comprehensive.

The overview presented in this report is thus a contribution to the discussion on remedies to the harms presented by generative AI, but extensive legal analysis will be necessary to determine the effect of these frameworks on generative AI.

## 3.1 Data protection law

The General Data Protection Regulation (GDPR)<sup>194</sup> applies to the processing<sup>195</sup> of personal data<sup>196</sup> by companies established in the European Union or by companies established outside the European Union, when these companies process the personal data of a data subject in the European Union (EU) or the European Economic Area (EEA).<sup>197</sup>

The obligations in the GDPR primarily apply to “controllers”, the entity that determines the purposes and means of processing personal data.<sup>198</sup> Some obligations are also laid on the “processors”,<sup>199</sup> an entity processing personal data on behalf of a controller. As mentioned in chapter 1.1.2, development and deployment of generative AI involves several actors at different stages of the process. It is crucial that the different actors in the generative AI actor chain clearly define their roles, to ensure compliance with the GDPR throughout the whole process.

As companies develop and deploy generative AI models, the GDPR could be applicable to at least three aspects of the system: the training data used to develop the generative AI model, the outputs from the generative AI model, and the generative AI model itself.

As described throughout this report, generative AI models analyse large amounts of data, typically scraped from the internet. Some of these data points are undeniably personal data, which means that the GDPR is applicable to the processing. Similarly, the GDPR applies to the processing of personal data from individuals’ prompts to generative AI models, associated with them through personal accounts or similar.

It is possible to use generative AI to generate images, text, videos, and audio relating to identifiable natural persons. The GDPR will therefore clearly be applicable to some output as well as the input.<sup>200</sup> This holds true

regardless of whether the generated information is correct, meaning that a deepfaked photo or an incorrect statement related to an identifiable individual is still personal data.

The way that a generative AI model works, the model will not necessarily include any personal data – there are no actual pictures of individuals in the model itself – but the output may be an identifiable image of a real person. However, even if the model does not contain personal data directly, researchers have been able to extract training data from large language models. Large language models are generally more vulnerable to such extraction than their smaller counterparts.<sup>201</sup> As mentioned above, the training data includes personal data as well, thereby potentially allowing for the extraction of personal data from the generative AI models. Some authors have argued that the possibility of extracting personal data from a model means that the model itself could be considered personal data.<sup>202</sup> Thus it is possible that the GDPR might apply to the models themselves, in addition to the input and output of the models.

Under the GDPR, the processing of personal data generally requires a legal basis.<sup>203</sup> There is a general prohibition on the processing of special category data, which includes categories such as personal data revealing racial or ethnic origin, political opinions, health, and biometric data.<sup>204</sup> In cases where the training data and/or output of a generative AI model includes special categories of personal data, the controller must have a legal basis that exempts this prohibition.

As of May 2023, some light has been shed on the legal bases that some developers claim for processing personal data for the development of generative AI models. After scrutiny from the Italian DPA,<sup>205</sup> OpenAI added a section in its privacy policy for international users, claiming legal bases such as performance of a contract and a broad legitimate interest to for example develop,

<sup>195</sup> Art. 4(2) GDPR.

<sup>196</sup> Art. 4(1) GDPR.

<sup>197</sup> Art. 3 GDPR.

<sup>198</sup> Art. 4(7) GDPR.

<sup>199</sup> Art. 4(8) GDPR.

<sup>200</sup> The GDPR does not apply to “natural persons’ processing of data in the

course of a purely personal or household activity”, cf. art. 3. Generation of output relating to a natural person as part of purely personal activity, where the output is not shared online, is therefore not necessarily regulated by the GDPR. The GDPR can still apply to the system owner.

<sup>204</sup> Art. 6 GDPR.

<sup>205</sup> Art. 9(2) GDPR.

improve or promote its services.<sup>206</sup> Google, on the other hand, has not yet released its chatbot Bard in the EU.<sup>207</sup> There has been speculation that this may be due to the GDPR,<sup>208</sup> and at the time of writing there is no mention of a legal basis for processing personal data in the context of Bard.<sup>209</sup>

In addition to needing a legal basis for processing personal data, there are various other relevant legal requirements for processing personal data for the training, development, deployment, and use of generative AI models. The discussion on the principles of data protection by design and by default,<sup>210</sup> data minimization<sup>211</sup> and purpose limitation<sup>212</sup> in the context of training machine learning models is nothing new, and the principles also apply to the training of generative AI models when personal data is involved.<sup>213</sup>

The principle of data minimization involves collecting and processing as little personal data as possible for the stated purposes of processing. Purpose limitation includes not using personal data for other purposes than stated at the point of collection, and to not store the personal data for longer than necessary to fulfil these purposes. As the training of generative AI models requires large amounts of data, and the models are often developed to be general purpose, these principles may come into conflict with the approach taken by many developers of generative AI models.

### 3.1.1 DATA SUBJECT RIGHTS

People whose personal data is processed (data subjects) have several rights under the GDPR. This includes the rights to erasure (having personal data deleted),<sup>214</sup> rectification (having personal data corrected),<sup>215</sup> and objection (protesting about the processing of personal data).<sup>216</sup>

It is still unclear how companies developing and deploying generative AI models will be able to fulfil requests to provide data subject rights in practice. After scrutiny of ChatGPT by the Italian Data Protection Authority, OpenAI introduced an opt-out mechanism to have personal data removed from the training data, and the possibility to correct inaccurate personal information. However, OpenAI clarifies in its privacy policy that “[g]iven the

technical complexity of how our models work, we may not be able to correct the inaccuracy”.<sup>217</sup> In other words, it is highly questionable whether it is technically feasible for OpenAI to provide data subject rights and comply with the GDPR.

It is also questionable whether an opt-out system such as the one implemented by OpenAI can be compliant with the GDPR. For the opt-out system to be effective, it would require the individual’s knowledge that a generative AI model was trained on their personal data. This is not evident to consumers unless they are frequent users of the generative AI models, and even then, it is unlikely that they would understand the extent of the processing of personal data.<sup>218</sup>

A significant hurdle that relates to deleting personal data from the training data is the sheer size of the data sets used to train generative AI models.<sup>219</sup> The work related to the collection, cleaning and preparation of data sets is generally not prioritized by AI practitioners, in favour of model development.<sup>220</sup> Consequently, companies’ ability to find and delete data traces of any individual is compromised by their lack of oversight and documentation of the data sets, which is at odds with data protection law.

<sup>210</sup> Art. 25 GDPR.

<sup>211</sup> Art. 5(1)(c) GDPR.

<sup>212</sup> Art. 5(1)(b) GDPR.

<sup>214</sup> Art. 17 GDPR.

<sup>215</sup> Art. 16 GDPR.

<sup>216</sup> Art. 21 GDPR.

<sup>218</sup> See however requirements of information for data subjects in art. 12-14 GDPR.

### 3.1.2 THE ITALIAN DPA'S DECISION CONCERNING CHATGPT

There have already been efforts to apply the GDPR to generative AI models. On March 31st 2023, the Italian Data Protection Authority (DPA) imposed a temporary limitation on OpenAI, the owner of ChatGPT, regarding the processing of personal data of Italian individuals. At the same time, the DPA opened an inquiry into the facts of the case.<sup>221</sup> This was due to several potential breaches of the GDPR, such as issues related to the processing of personal data about end users of the ChatGPT service, processing of personal data in relation to the training of the model, as well as processing of personal data during content generation. As a result, OpenAI temporarily blocked access to ChatGPT for individuals located in Italy. While this addressed some of the data protection issues described by the DPA, it was for example still possible for peoples outside of Italy to generate personal data concerning Italian citizens.

Some of the potential breaches pointed out by the Italian DPA have more far-reaching consequences than others. It is possible for OpenAI to implement measures to address issues such as data breaches and age verification mechanisms without altering its model significantly. Generation of inaccurate personal data seems more difficult to address, even though it is in line with OpenAI's general attempts to increase accuracy in its models. As mentioned above, OpenAI already states that it may not be able to correct inaccuracies,<sup>222</sup> and it seems very unlikely that the company would be able to guarantee accurate personal data about individuals, either intrinsically in their model or through content moderation, as discussed above and in chapter 2.3.2.

The final and most damning issue that the Italian DPA raised is that OpenAI appeared to have no legal basis to process personal data about Italian citizens to train its model. As the GDPR is harmonized in the EU, this would effectively mean that OpenAI did not have a legal basis

to train its generative AI model on personal data from any data subjects in the EU or the EEA. While OpenAI has not shared information about its training data, it is safe to assume that it contains personal data about data subjects in the EU and the EEA, for instance scraped from the internet.

It is technically possible to prepare new data sets for the training of subsequent GPT models, cleaning the data sets to remove personal data about data subjects in the EU and the EEA, but this would be extremely time and resource intensive, and probably halt development significantly. In any case, this issue raises the question of whether general purpose generative AI models<sup>223</sup> and the GDPR can coexist in their current form.

On April 28th 2023, ChatGPT was reinstated in Italy, after OpenAI introduced various data protection measures, such as the opt-out mechanism described above, a mechanism to exercise the right to erasure of personal data, a new information notice including the legal bases used for processing, and age specification requirements.<sup>224</sup> While OpenAI has on paper introduced additional data protection measures, it is unclear how consumers would be able to make effective use of their rights, such as the right to opt-out from their personal data being used to train the generative AI model, in practice.

Even if OpenAI were to have legitimate interests that are overridden by the fundamental rights of data subjects,<sup>225</sup> OpenAI seemingly did not conduct this balancing act before releasing its generative AI model to the public. OpenAI's compliance with for example the principles of lawfulness of processing or accountability<sup>226</sup> therefore appears dubious at best. The Italian DPA's apparent acceptance of these claims could prove problematic, as compliance with the GDPR goes beyond the quick fixes deployed by OpenAI.

<sup>223</sup> Umbrella term for AI systems designed to perform a wide range of tasks across different domains.

<sup>225</sup> Art. 6(1)(f) GDPR.

<sup>226</sup> Art. 6 and 5(2) GDPR.

It seems evident that the GDPR's protections of individuals will require extensive analysis and swift enforcement to be effective. Even if the GDPR can serve to protect individuals, the regulation has been criticized for being slow and complicated to enforce particularly in cross-border cases.<sup>227</sup> The European Data Protection Board (EDPB) has however established an EU-wide "task force" to coordinate investigations and enforcement on ChatGPT, so there will evidently be more legal development.<sup>228</sup> The French DPA has received several complaints and published an action plan involving ChatGPT.<sup>229</sup> Both German and Spanish DPAs are also considering actions.<sup>230</sup>

Although OpenAI and ChatGPT has been the main focus of GDPR enforcement in the field so far, others are likely to follow as different models become widely used. There is clearly appetite for enforcement, and so the GDPR will be an important legal framework in ensuring generative AI that respect data protection rights for all data subjects.

## 3.2 Consumer law

The Unfair Commercial Practices Directive (UCPD)<sup>231</sup> lays out the legal provisions governing business-to-consumer practices in the EU and EEA. It is technology-neutral, applies to all business-to-consumer transactions, and is meant to work as a catch-all to protect the consumer's decision-making in the commercial environment. The directive targets traders' practices, to ensure these commercial practices<sup>232</sup> are not unfair, often through requirements on disclosure and openness.

Certain commercial practices are banned outright, through the blacklist of unfair commercial practices in the UCPD Annex 1. In addition to the banned practices of Annex 1, there are several broad, discretionary legal provisions in the UCPD. A commercial practice is unfair if it is misleading<sup>234</sup> or aggressive,<sup>235</sup> causing (or likely to cause) an average consumer to take a transactional decision she would not have taken otherwise.

Transactional decisions have been broadly defined and include consumer decisions such as adding an item to a virtual basket or entering a shop. In the most recent guidelines to the UCPD, published in December 2021, the EU Commission also included examples such as continuing to use a service by browsing or scrolling,<sup>236</sup> thus apparently broadening the scope of the transactional decision test of the UCPD to entail business practices at the core of the attention economy. As the Guidelines are not legally binding, it is not yet clear how they will be interpreted by Consumer Authorities and courts in practice.

There is also a general clause under which a practice is unfair if it is contrary to the requirements of professional diligence and it distorts, or is likely to distort, the average consumer's economic behaviour. This serves as a safety net to ensure that unfair practices not covered by the blacklist, misleading practices or aggressive practices can still be subjected to the fairness assessment of the UCPD.<sup>237</sup> The requirement of professional due diligence could also serve as a bridge to other legal frameworks, such as non-discrimination law, making it possible to integrate jurisprudence from such laws in the consumer law context.<sup>238</sup>

Generative AI as stand-alone models or embedded in other, consumer-facing services can potentially be addressed by the UCPD in several ways, either argued in a traditional monetary sense, or by keeping consumers on the service. In any case, the applicability of the UCPD depends on the generative AI model being used in the context of a commercial practice.

Bing is currently employing advertisements in their generative AI search,<sup>239</sup> which requires distinct labelling to ensure that it is not a misleading practice. If a text generator is used to persuade the consumer to stay engaged with the service, for instance through persistent communications, particularly targeting the consumer's identified weaknesses (as may be the case for chatbots programmed to generate and simulate romantic involvement), this could also amount to an aggressive practice. As mentioned in chapter 2.1.4.3, companies are already

<sup>232</sup> Art. 2(b) UCPD.

<sup>233</sup> Art. 4(d) UCPD.

<sup>234</sup> Art. 6-7 UCPD.

<sup>235</sup> Art. 8-9 UCPD.



attempting to integrate generative AI into their shopping experiences, which could potentially mislead consumers into buying a product by providing inaccurate information.

There have been calls to address consumer challenges stemming from generative AI through the UCPD. The European consumer organization BEUC<sup>240</sup> addressed a letter about generative AI, and particularly text generators, to DG JUST and the Consumer Protection Cooperation Network (CPC Network) on April 21st, 2023.<sup>241</sup> In the letter, BEUC brought attention to various ways in which generative AI is deployed that may influence consumer behaviour in a way that is in breach of the UCPD, also taking account of vulnerable groups such as children.

There is clearly a need for consumer authorities to consider and address illegitimate use of generative AI in commercial contexts as they are employed in an increasing number of services. The UCPD may be used to tackle certain challenges related to generative AI in the context of commercial practices. Particularly relevant examples may be if a generative AI model is used in a way that provides consumers with false or misleading information about for example products, or if the system owner omits information from the terms of services.

At the same time, there are some potential limitations to the UCPD's applicability to generative AI. Presumably, generative AI may often be used to increase consumer interaction with a trader's service and boost engagement by capitalising on the AI's capability to seem like a sentient interlocutor despite its obvious limitations. This may mislead consumers to dedicate much more time and attention to such services than they normally would, absent the illusion created by the model, for example if a chatbot is designed to simulate romantic emotions toward the consumer. However, the UCPD's usefulness against practices unfairly extracting consumer attention and engagement is not yet certain. This would require a broad understanding of what constitutes "transactional decisions" that so far is only based on the European Commission's (non-binding) Guidelines, rather than the letter of the law. As such, many relevant practices are not clearly covered by the UCPD.<sup>242</sup>

There is also the question of whether the remedial mechanisms of the UCPD are sufficient, unless they are applied in a way to also require a 'fairness by design' to

be offered by digital services as a necessary pillar of professional diligence, rather than focusing on disclosure requirements. The EU Commission should take advantage of the ongoing digital fitness check, which is an evaluation of the current EU consumer law's fitness to ensure a high level of protection,<sup>243</sup> to address these problems.

### 3.2.1 CONSUMER LAW IN A U.S. SETTING

In the United States, the Federal Trade Commission (FTC) enforces the FTC Act, which gives the FTC the authority to enforce violations of unfair and deceptive trade practices.<sup>244</sup> Unlike the European Consumer Authorities, the FTC is mandated to prescribe rules to target specific deceptive practices or unfair methods of competition within the scope of the FTC Act.<sup>245</sup> As of June 2023, the FTC is in the midst of a broad rulemaking process that can create new rules proscribing specific practices that are *always* unfair and deceptive. This means that the FTC could react more forcefully and specifically to new technologies on the market.

The FTC has already issued guidelines on AI in 2021, calling for transparency, unbiased results, accountability, and more.<sup>246</sup> In 2023, amidst the surge of generative AI products and services, the FTC also issued a statement to businesses, reminding them of the need to advertise AI in a truthful and responsible way.<sup>247</sup>

Although the FTC has not yet undertaken enforcement actions against companies deploying or training generative AI, it has had cases with strong equitable remedies including Algorithmic Disgorgement.<sup>248</sup> The remedy requires the company to delete data and models/algorithms built on that data when the consumers' rights were violated in the collection process. This powerful remedy may improve the practices of companies that fear being forced to delete their models.

The Center for AI and Digital Policy filed a complaint to the FTC on March 3rd, 2023, asking for a moratorium on the release of further commercial versions of ChatGPT beyond GPT4, and for rulemaking relating specifically to generative AI.<sup>249</sup> The FTC does not open cases for individual consumers, but can base investigations on such complaints and reports. It therefore seems likely that the FTC might react to this particular complaint due to the mass adoption of generative AI, in order to minimize consumer harm.

<sup>240</sup> A European consumer umbrella organization of which the Norwegian Consumer Council is a member.

On the 25th of April 2023, the FTC, Consumer Financial Protection Bureau (CFPB), Equal Employment Opportunity Commission (EEOC), and Department of Justice's Civil Rights Division announced that they intend to enforce against "discrimination and bias in automated systems,"<sup>250</sup>

while the White House announced several initiatives on May 5th, including committing the public sector to mitigating risks.<sup>251</sup>

## 3.3 General product safety law

Product safety legislation is meant to ensure that products placed on the market are safe for consumers to use. Current European general product safety law is based on the General Product Safety Directive (GPSD).<sup>252</sup> By the end of 2024, the General Product Safety Regulation (GPSR)<sup>253</sup> will replace the GPSD. Both legal instruments are relevant in the context of generative AI.

### 3.3.1 THE GENERAL PRODUCT SAFETY DIRECTIVE

The GPSD complements sector specific legislation and applies to any risks from a product not covered by other laws.<sup>254</sup> In practice, the GPSD performs the role of a safety net, ensuring safety requirements for all products on the European market.

The legislation requires that producers place only safe products on the market.<sup>255</sup> While the Directive's definition of a product is broad enough to theoretically cover harms resulting from software linked to a product,<sup>256</sup> the scope does not explicitly include or exclude software. Its applicability to GPT-models and other purely software based generative AI models is therefore still uncertain.

A product is deemed safe when it does not present any risk, or more than minimal risk to consumers' safety and health under normal and foreseeable use.<sup>257</sup> This has traditionally covered physical impacts on persons, such as physical injuries or property damage, and mental health is not explicitly mentioned in the GPSD. While some argue that intrinsic mental health risks from products could be covered by the GPSD,<sup>258</sup> the lack of explicit reference makes its applicability to mental health risks more uncertain.

Competent authorities are required to consider whether products are in fact safe, even after they have been

placed on the market.<sup>259</sup> Such a consideration must take account of the precautionary principle, meaning that a product may be presumed unsafe in the absence of scientific certainty about potential harms and harmful effects of the product.

As explained throughout this report, it seems clear that generative AI may in fact pose considerable risks to consumers, especially mental health risks. Such risks may for example arise from the generation and subsequent dissemination of inaccurate personal data or deepfakes, from the deployment of highly manipulative and personalized generative AI models, or in situations where consumers use generative AI models for mental health or medical advice purposes.

There have been moves to apply the GPSD to generative AI, through alerting safety authorities to investigate safety risks of generative AI. The European consumer organization BEUC sent a letter to the Consumer Safety Network on April 12th, 2023 in this regard.<sup>260</sup> The letter particularly brought attention to risks to consumer mental health.

### 3.3.2 THE GENERAL PRODUCT SAFETY REGULATION

A new General Product Safety Regulation (GPSR) has been approved by the EU, which will come into force by the end of 2024. The new regulation will repeal the GPSD and widens the scope of products by bringing software into the scope, as well as explicitly mentioning mental health.<sup>261</sup> It will also require any producer to consider a product's evolving, learning and predictive functionalities when assessing the risks of the product,<sup>262</sup> which is clearly relevant in the context of generative AI products.

<sup>254</sup> Art. 1(2) GPSD.

<sup>255</sup> Art. 3(1) GPSD.

<sup>257</sup> Art. 2(1)(b) GPSD.

<sup>259</sup> Art. 8(2) GPSD.

<sup>261</sup> Cf. recital 19 GPSR.

<sup>262</sup> Art. 6(1)(h) GPSR.

The applicability of the GPSD to generative AI models may be uncertain, but it seems clear that the GPSR will apply to generative AI models. It is necessary that safety authorities take preventive action to address harms stemming from generative AI to the degree possible

under the current legal framework. This will also help safety authorities prepare thoroughly to enforce the GPSR as soon as it comes into force and starts applying to generative AI models.

## 3.4 Competition law

At its core, EU competition law serves to prevent anti-competitive practices so markets remain contestable, and consumers can benefit from lower prices, better quality of products and services, and more choice and innovation.

The essence of EU competition law is found in the Treaty on the Functioning of the European Union (TFEU), although implemented through other regulations. First, it is prohibited for companies to make anti-competitive agreements.<sup>263</sup> Second, companies may not abuse their dominant position.<sup>264</sup>

The concept of a “dominant position” is highly context-based and will depend on how the “relevant market” is defined. This includes, for example, the availability of alternative products and consumers’ willingness to switch to these alternative products.<sup>265</sup> As discussed in chapter 2.1.3 above, the deployment of generative AI carries with it a risk of power concentration in the hands of a few actors. This may lead to certain companies becoming dominant within their respective markets, for example generative AI based search engines, shopping assistants, etc.

The digital sector notoriously consists of a very few, massive actors, often referred to as big tech. It is crucial that any emerging market relating to generative AI faces scrutiny from antitrust agencies early, to avoid a similar emergence of highly dominant companies in this market, who may be tempted to abuse their position. It is noteworthy that several big tech companies are investing heavily in generative AI.

Competition enforcement agencies clearly have a role to play to address some of the harms outlined in this report. The UK Competition and Markets Authority (CMA) has opened an initial review of competition and consumer protection considerations in the development and use of generative AI. Other competition authorities should closely monitor future developments in this sector and be ready to intervene early if they notice behaviours or practices falling within their remit, which could be anti-competitive. This would help ensure that the emerging generative AI sector remains fair and competitive.

## 3.5 Content moderation

The Digital Services Act (DSA)<sup>266</sup> is a new EU regulation, aiming to improve the mechanisms for removing illegal content, and to protect individuals’ rights, including freedom of speech and a high level of consumer protection. It will be an important tool in introducing new content moderation to online services.

The DSA applies to online intermediary services, meaning conduit services, caching services and hosting services.<sup>267</sup> In practice, this means that the DSA applies to services

connecting consumers to goods, services, and content – such as online marketplaces, social media platforms, cloud hosting services, and internet access providers.

It will be fully applicable to all entities in its scope in February 2024, while very large online platforms (VLOPs) and very large online search engines (VLOSEs), some of which have already been designated by the European Commission,<sup>268</sup> will have to oblige by the new rules as of the end of the summer 2023.

<sup>263</sup> Art. 101 TFEU.

<sup>264</sup> Art. 102 TFEU.

<sup>267</sup> Art. 3(1)(g) DSA.

Generative AI models are not clearly covered by the DSA. The most relevant type of service covered by the DSA is “hosting services”.<sup>269</sup> Even then the applicability of the DSA is not clear, since the content provided by generative AI models is largely generated by the models themselves, rather than consumers or other third parties.

While generative AI models as stand-alone services may not be covered by the DSA, the DSA may be applicable to companies that wish to embed generative AI models in their platforms and services. The integration of ChatGPT

in the search engine Bing, which is a designated VLOSE, may for example in effect trigger the content moderation requirements of DSA for the generated content.

Any content generated by generative AI and subsequently shared or stored by consumers on services covered by the DSA will similarly be covered by content moderation requirements. In both cases, the DSA appears to apply to the downstream provision and utilization of the generated content, rather than the content generation model itself.

## 3.6 The draft Artificial Intelligence Act

In April 2021, the EU Commission published a proposal for the Artificial Intelligence Act (AIA), laying down harmonized rules across the EU and the EEA “to foster the development, use and uptake of artificial intelligence in the internal market”.<sup>270</sup>

An EU legal framework aimed at regulating AI should be expected to also regulate generative AI. However, the Commission’s proposal for the AIA was published before the widespread adoption of generative AI during the winter 2022/2023. In the aftermath, discussions among EU lawmakers have focused heavily on how to properly regulate generative AI as part of the AIA.

As the AIA is not yet complete, it is not yet certain how it will apply to generative AI in practice. In the following, the Commission’s draft, as well as the Council and Parliament positions to the AIA as of May 2023 are briefly outlined. The overview begins with relevant aspects from the Commission’s proposal for the AIA.

### 3.6.1 THE EU COMMISSION’S PROPOSAL

The Commission’s proposal for the AIA (hereafter “Draft AIA”) applies to any provider placing an AI system on the market.<sup>271</sup> AI systems are defined broadly, as systems

generating output based on machine learning approaches, logic- and knowledge-based approaches, or statistical approaches.<sup>272</sup> In other words, the scope of the Draft AIA is broad, encompassing many types of systems.

The AIA has a risk-based approach, regulating different types of AI systems based on their risks to individuals or society. Certain exclusively listed practices are prohibited,<sup>273</sup> and may never be put on the European market. AI systems may also be classified as high-risk, for example if they are among the high-risk systems listed in Annex 3.<sup>274</sup>

Most of the Draft AIA focuses on regulating high-risk AI systems and setting out legal requirements for AI operators of these systems.<sup>275</sup> This includes legal requirements such as creating a quality management system,<sup>276</sup> including a risk management system,<sup>277</sup> to meeting data quality criteria,<sup>278</sup> accuracy, robustness, and cybersecurity measures,<sup>279</sup> as well as creating technical documentation.<sup>280</sup>

Notably, any provider of a system that falls outside the scope of high-risk systems have few, if any, requirements laid upon them through the Draft AIA. There are some limited transparency requirements for applications such as chatbots and for deepfake material.<sup>281</sup> All providers of

<sup>271</sup> Art. 2(1)(a) Draft AIA.

<sup>272</sup> Art. 3(1) Draft AIA, cf. Annex 1.

<sup>273</sup> Art. 5 Draft AIA.

<sup>274</sup> Art. 6 Draft AIA.

<sup>275</sup> See Draft AIA Title III.

<sup>276</sup> Art. 17 Draft AIA.

<sup>277</sup> Art. 9 Draft AIA.

<sup>278</sup> Art. 10 Draft AIA.

<sup>279</sup> Art. 15 Draft AIA.

<sup>280</sup> Art. 11 Draft AIA.

<sup>281</sup> Art. 52 Draft AIA.

systems that are not high-risk may also voluntarily abide by the requirements for high-risk systems,<sup>282</sup> but have no legal obligation to do so. Thus, the Draft AIA targets a very broad range of AI systems, while imposing obligations on very few of them, and precluding Member states from imposing additional obligations.<sup>283</sup> This makes the scope of high-risk AI systems especially important. Additionally, the Draft AIA contains very limited rights for consumers.

It is unclear how generative AI systems fit into the Commission's Draft AIA from 2021. It is plausible that they would have fallen inside the scope of the Draft AIA generally. For more specific requirements, the generative AI systems would have to be related to one of the high-risk categories in Annex III, be used in the context of chatbots or deepfakes for limited transparency requirements or be related to a prohibited practice.<sup>284</sup>

### 3.6.2 THE EU COUNCIL'S POSITION ON THE AIA

In the Council's position on the AI Act (hereafter, "Council Position"),<sup>285</sup> general purpose AI is addressed. The definition reads as follows:

*an AI system that - irrespective of how it is placed on the market or put into service, including as open source software - is intended by the provider to perform generally applicable functions such as image and speech recognition, audio and video generation, pattern detection, question answering, translation and others; a general purpose AI system may be used in a plurality of contexts and be integrated in a plurality of other AI systems;*<sup>286</sup>

This would apply to all types of generative artificial intelligence touched upon in this report.

In the Council Position, general purpose AI would be subject to high-risk obligations if they "**may** be used as high-risk AI systems or as components of high risk AI systems".<sup>287</sup> Whenever the AI provider explicitly excludes all high-risk uses in the instructions or information relating to the generative AI, the general purpose AI system is exempted.<sup>288</sup> This exemption may only apply when the

exclusion is made in good faith.<sup>289</sup>

In practice, it would be very hard for providers to ensure that their system can never be used in high-risk settings, as defined in the Council Position Annex 3. The extent of the "good faith" requirement is therefore crucial; either it practically requires that all generative AI systems are subject to high-risk obligations, or it may prove to be too low of a threshold, serving as the basis for meaningless disclaimers from developers. In any case, the Council Position's effect on generative AI systems is uncertain.

### 3.6.3 THE EU PARLIAMENT'S POSITION ON THE AIA

As of May 2023, the EU Parliament's position is still being negotiated. The LIBE and IMCO committees of the EU parliament approved a compromise position on May 11th.<sup>290</sup> It will be voted over in plenary in mid-June. The following therefore constitutes the parliament's presumed position at the time of writing (hereafter "Parliament Position").

Overall, the European Parliament's position significantly improved the European Commission's proposal. Consumers are granted new rights including a right to be informed when being subject to a decision from a high-risk AI system,<sup>291</sup> a right to complain to an authority about an AI system,<sup>292</sup> and the right to bring a supervisory authority to court if it fails to take action.<sup>293</sup> Consumers were also given the right to ask for collective redress when an AI system has caused harm to a group of consumers.<sup>294</sup>

When it comes to generative AI, rather than focusing on general-purpose AI like the Council's position, the European Parliament introduces a new concept: 'foundation model'. The Parliament Position defines 'foundation model' as an AI model that is "trained on broad data at scale, is designed for generality of output, and can be adapted to a wide range of distinctive tasks".<sup>295</sup> All providers of foundation models have additional obligations, regardless of whether the models are provided as a standalone model or embedded in a system, or whether the models are open source or closed source.<sup>296</sup> These obligations include requirements to identify and reduce risks to for example health, safety and the rule of law, data governance measures, appropriate levels of for example

<sup>282</sup> Art. 69 Draft AIA.

<sup>284</sup> Art. 5 Draft AIA.

<sup>286</sup> Art. 3(1)(b) Council Position.

<sup>287</sup> Art. 4b(1) Council Position.

<sup>288</sup> Art. 4c(1) Council Position.

<sup>289</sup> Art. 4c(2) Council Position.

<sup>291</sup> Art. 68c Parliament Position.

<sup>292</sup> Art. 68a Parliament Position.

<sup>293</sup> Art. 68b Parliament Position.

<sup>294</sup> Art. 68d Parliament Position.

<sup>295</sup> Art. 3(1c) Parliament Position.

<sup>296</sup> Art. 28b Parliament Position.

performance, predictability, interpretability throughout its life cycle, energy efficiency measures, and technical documentation.<sup>297</sup>

Models used as a basis for generative AI systems are clearly meant to be encompassed by the definition of foundation model. The Parliament Position refers explicitly to such systems in the article placing obligations on providers of foundation models. Foundation models used in generative AI systems are placed under additional obligations of transparency, adequate safeguards against the generation of unlawful content, and the publication of a “sufficiently detailed summary of the use of training data protected under copyright law”.<sup>298</sup>

#### 3.6.4 THE ARTIFICIAL INTELLIGENCE ACT MUST PROTECT CONSUMERS

Since generative AI models are not built for a particular context, and they allow for wide-scale use, some authors have argued that generative AI models do not fit well into the risk-based system of the Draft AIA.<sup>299</sup> Instead, they argue targeted measures such as system risk monitoring should be considered. At the same time, as outlined throughout this report, generative AI systems pose significant risks that must be mitigated at the development stage of the system, rather than at the time the system is placed on the market, or after the system is placed on the market.<sup>300</sup>

It is not yet certain how the AIA will apply to generative AI. With the AIA, however, European lawmakers have a unique chance to introduce enforceable guardrails to protect consumers against the risks of generative AI. That chance must be used effectively in the months to

come before the AIA is finalized and should address the harms outlined in this report both through introducing consumer rights, and through obligations on the entire generative AI actor chain.

EU lawmakers must ensure that industry lobbying does not water out obligations and rights for consumers in the AIA. According to a report from Corporate Europe Observatory, industry lobbying efforts significantly weakened several relevant provisions in the proposed regulation from the Commission, including pushing for excluding general purpose AI systems from the regulation.<sup>301</sup> EU lawmakers must be vigilant to avoid falling for lobbying tactics, that will no doubt increase during the final stage of negotiations.

While lawmakers finalize the AIA, other key enforcement agencies must also ensure the safety and rights of consumers. The AIA will not be fully applicable for several years,<sup>302</sup> and in the meantime, it is necessary for enforcement agencies of other legal frameworks to protect consumers from the harms of generative AI, as outlined above.

**“With the AIA, however, European lawmakers have a unique chance to introduce enforceable guardrails to protect consumers against the risks of generative AI.”**

## 3.7 Liability

There are several relevant liability laws in the EU, meant to ensure that consumers receive fair compensation when defective products lead to harm. Some legal instruments are already in force, while others are still being negotiated.

### 3.7.1 PRODUCT LIABILITY DIRECTIVE

Product liability rules enable consumers to claim compensation for damage caused by a defective product. The current EU liability rules – the Product Liability Directive (PLD) – was adopted in 1985 and it is not clear whether it applies to generative AI.

<sup>297</sup> For the complete list of requirements, see art. 28b(2) Parliament Position.

<sup>298</sup> Art. 28b(4).

<sup>302</sup> The exact timing differs between the different positions of the EU institutions. At the earliest the AIA will be fully applicable after 24 months after entering into force. This effectively means it may not be fully

applicable until at the earliest April/May 2026, if there is a trilogue agreement by January 2024.

First, there is no consensus on whether the PLD applies to digital services and software such as generative AI.

Secondly, even if the PLD applied to digital services such as generative AI, a court ruling from the Court of Justice of the European Union established that information provided by a product is not covered by the PLD.<sup>303</sup> Since the output of generative AI essentially is information in the form of speech, text or pictures, the fact that information is not covered by the PLD means that generative AI is most likely not covered by the PLD.

### 3.7.2 REVISED PRODUCT LIABILITY DIRECTIVE

The European Commission has introduced a proposal for an updated product liability directive (revised PLD). The updated directive is also meant to cover software, including AI systems.<sup>304</sup>

The revised PLD proposal, similarly to the PLD, operates with a non-fault-based liability scheme. While consumers will not have to prove fault with the operator or the producer of a product, consumers will have to prove the relevant defect in a product, the consumer harm and the causal link between the defect and the harm.

It remains to be seen how the recent ruling from the Court of Justice on the current PLD will apply in the context of the revised PLD, thereby whether information provided by a product will be covered by the revised PLD or not.

In any case, for generative AI systems, most of the potential consumer harms are non-material, as outlined in chapter 2. Such harms are exempt from the PLD, which only covers material damage.

In total, neither the current PLD nor the revised PLD appear well suited to allow consumers compensation for harms from generative AI systems. The final wording of the PLD proposal will however have to be considered.

### 3.7.3 AI LIABILITY DIRECTIVE

In parallel with the AIA, and as an addition to the PLD, the European Commission has proposed the AI Liability Directive (AILD), which is meant to provide consumers with the possibility to claim compensation for harms caused by AI systems. While the EU Commission has proposed a draft, it is likely that the Directive will not be finalized until the AIA has been adopted.

The AILD proposal gives consumers the possibility to claim compensation for all material harms, and non-material harms if allowed within national legal frameworks. However, there are serious limitations to the proposal which will significantly reduce its effectiveness in providing compensation for consumer harms.<sup>305</sup>

Consumers that wish to claim compensation for harms caused by AI systems are required to prove fault of the AI system operator. Proving fault in the context of the AILD means that the consumer will have to prove that the AI system operator is not operating in accordance with the EU rules, including the AIA. Proving such non-compliance will require high technical and legal knowledge, none of which regular consumers have or should be expected to have. Since proving fault is a prerequisite for other mechanisms in the AILD, such as a presumption of causality between fault and the output leading to harm, this limitation is substantial. To make the AILD effectively protect consumers, non-fault-based liability should be established in claims from consumers and the burden of proof should be reversed.<sup>306</sup>

When it comes to generative AI, it is still uncertain whether the AIA will classify generative AI as an 'AI system'. If it does, since the definition of AI systems in the AILD refers to the definition of the AIA, then the AILD would be applicable to generative AI. However, compensation claims for damages created due to inaccurate information (in the form of text, images, or audio) are not harmonised under the AILD proposal. This means that compensation claims for consumers related to generative AI would therefore need to be assessed at the national level on a case-by-case basis.

The AILD is still early in the political process, and EU lawmakers have room to amend the proposal in a manner that would give consumers real options to seek compensation from harms arising from generative AI, regardless of national rules. This is necessary to increase consumer protection in the face of the harms outlined in chapter 2.

<sup>303</sup> Case C-65/20, VI v Krone (2021), <https://curia.europa.eu/juris/document/document.jsf?jsessionid=7A0662FAD49ED462BE89A81594FAF809?text=&docid=242561&pageIndex=0&doclang=EN>.

## 3.8 Industry standards and guidelines

Industry players are already developing guidelines that are meant to increase the transparency of both development and use of generative AI models.<sup>307</sup> There have also been industry calls to halt the development of new generative AI models.<sup>308</sup> The calls for a halt in development typically focus on the risks of very advanced models, and has coincided with a voluntary halt of development of GPT5 from Open AI.<sup>309</sup> The many risks from current generative AI models such as systems based on GPT4, as identified and discussed in chapter 2 of this report, are however addressed insufficiently.

There are increasingly also calls for creating voluntary codes of conduct for developers and deployers of generative AI.<sup>310</sup> For the EU in particular, the Commission is aiming to create a pact with companies ahead of the new rules in the AI Act. Policy makers in the EU are reportedly planning to co-create a code of conduct “within months”, which means the codes of conduct would be created or negotiated at the same time as the trilogue negotiations of the AIA. Industry representatives, such as Google, would in this way be perfectly placed to increase their lobbying efforts.

The envisioned process creates a dual risk. First, the European Commission has a role to play in trilogue negotiations and it cannot fulfil this role impartially if, at the same time, it is negotiating a code of conduct or

other self-regulatory rules with the industry and third countries on the same topic. Secondly, it is unclear what requirements a voluntary agreement can include when the legal requirements for these actors in the EU are not yet defined. There is an obvious risk that the voluntary commitments will not be in line with the final legal text. The voluntary codes of conduct would be strongly affected by industry players’ views on feasibility and possibility to monetize products, instead of consumer and human rights. Finally, the AIA may also be unduly affected by industry players. This is unacceptable and must not be allowed.

Industry standards and guidelines are unsuited to tackling the risks stemming from the deployment of GPT-models already on the market, as they tend to act as a lowest common denominator, lacks sufficient enforcement mechanisms and independent oversight. Instead of relying on voluntary industry commitments while the AIA is not yet applicable, authorities should focus on the enforcement of existing laws such as consumer protection, data protection or product safety legislation. Policymakers and lawmakers should on their part strive to avoid self-regulatory regimes.





# **4. THE WAY FORWARD**

Throughout this report, we have outlined significant harms and challenges related to the development, training, deployment, and use of generative AI. These are not hypothetical risk of future dystopias, but tangible harms that affect people and populations today.

We believe that while these issues are concerning, they are not insurmountable. Many of the problems regarding generative AI are echoes of well-known issues from other sectors, but the rapid development and adoption of generative AI models mean that it is pertinent to take actions to address the harms. We cannot afford to wait until the technology is so embedded in our lives and social structures that it is too late to change the direction of its development and use.

Technology is not an untameable beast but must be adapted to and shaped by the rules and values of democratic societies. To ensure that generative AI is developed and used in accordance with consumer and human rights, it is clearly insufficient to rely on companies to regulate themselves. It is the responsibility of policymakers and enforcement agencies to set boundaries for how the technology is trained, developed, deployed, and used. Therefore, lawmakers must not pass legislation based on what industry actors state is technically feasible, but rather what is necessary to provide safe and consumer-centric technology in the years to come.

Below we set forth some fundamental principles that we believe should be at the heart of how society approaches generative AI. This is followed by several action points for enforcement agencies, policymakers, and lawmakers, and enforcement agencies. We hope that these points will provide a blueprint for a human-centric approach to the technology.

## 4.1 Consumer rights principles that are key for safe and responsible AI

To ensure that generative AI is safe, trustworthy, fair, equitable, and accountable, there is a need for overarching principles that address consumer rights. The principles set out below provide a foundation for how policymakers and enforcement agencies should approach the opportunities and pitfalls of generative AI.

Many of the principles are already defined in current consumer law, but we urge policymakers and enforcement agencies to ensure that they are in fact the foundation for development and deployment of generative AI. This is key to ensure a technological landscape that respects fundamental consumer rights today and in the coming years.

- **Consumer rights must be respected.** The onset of generative AI must not undermine or displace already established consumer- and human rights, such as a right to information and transparency, fairness and non-discrimination, safety and security, privacy and personal data protection, and redress.
- Consumers must have the **right to object** and **to an explanation** whenever a generative AI model is used to make decisions that have a significant effect on the consumer.
- Consumers must have a **'right to be forgotten'** to have personal data **deleted** from generative AI models, to **rectify** harms from for example having false information produced about them.
- Consumers must have a **right to interact with a human instead of generative AI** where this is relevant, for example in customer service contexts. This should not incur additional costs on the consumer, so that consumers are not treated disparately and unfairly based on their ability to pay.
- Consumers must have a **right to redress and compensation for any damages** suffered from the use of generative AI.
- Consumers must have a **right to collective redress**, and to be represented by consumer organizations and other civil society groups with exercising their rights.
- Consumers must have a **right to complain to supervisory authorities or launch legal actions** in court when use of a generative AI model is in breach of the law.
- Developers and deployers of generative AI models **must establish systems to ensure that these rights are available** to consumers in practice.

## 4.2 Policy recommendations

The decision on how to integrate technology into society is an inherently political question. Elected officials and governments have a responsibility to ensure that technology serves the people, rather than the whims of a small number of companies. A consumer-oriented technology policy means that people and societies should not be used as testing laboratories for experimental technologies. The lessons learned from the broader transition into a digital society without sufficiently accounting for the rights of citizens and the effects on society must inform how governments approach generative AI.

To ensure responsible, fair, and accountable innovation on the terms of society, we need robust policies that are future proof, rather than being swept up in the hype and then having to course correct in the aftermath. Below, we present several action points for how governments and policymakers should approach generative AI and similar technologies.

### 4.2.1 CALLS TO ACTION AND EMPOWERMENT OF ENFORCEMENT AGENCIES

While emerging technologies such as generative AI are sometimes described as a regulatory wild west, there are already comprehensive legal frameworks in place. We believe that many of these regulations are already suited to address a number of the issues described in chapter 2 of this report. However, to effectively protect people from exploitation, discrimination, and other abuses of power, these laws must be enforced.

Effective enforcement requires that enforcement agencies have the necessary powers, expertise, and resources, and that entities that are unable to or refuse to comply are sufficiently dealt with. In this section we present several necessary approaches and prerequisites for enforcement agencies to use their existing tools to shape technology in a consumer-friendly way.

- **Enforcement agencies must not wait for upcoming regulation.** Instead, they must **immediately investigate** generative AI systems **and apply relevant legal provisions** from their respective legal frameworks, such as data protection, competition, product safety and consumer law.
- **Collaborative cross-sectorial investigations**, where several enforcement agencies are involved in the same investigation, may be necessary to manage the risks stemming from generative AI. It may be necessary to appoint a coordinator for algorithmic enforcement, to ensure progress in such collaborative investigations.
- **Enforcement agencies should be empowered to conduct post-market surveillance of** generative AI models and have the option to order **product recalls** or the **deletion of algorithmic systems** or parts thereof that fail to comply with relevant legislation. Such orders should be accompanied by significant monetary fines to deter bad practices.
- **Enforcement agencies must have all necessary resources to enforce infringements** of their respective legal frameworks, including personal and technical competence and the necessary technical tools. With the deluge of AI-generated content, scaling up market surveillance and enforcement will be necessary.
- **Transnational and national technological expert groups should be established to support enforcement agencies** in enforcement endeavours.
- **Research** must be conducted on how to **augment enforcement by using technology.**

#### 4.2.2 DECISION MAKERS – STRATEGIC MEASURES

- **Governments must take critical perspectives on generative AI into account in their national AI strategies.** Over-arching principles to promote safe and human-centric generative AI must be built in from start rather than as add-ons, because add-ons are costly, and mistakes will erode trust.
- **Governments must adopt a critical and precautionary approach to using generative AI in the public sector.** The public sector has a particular responsibility to employ generative AI in a legal and trustworthy manner, and public procurement should be utilized to actively influence providers of stand-alone generative AI software or systems with embedded generative AI. In particular, the public sector should require transparency, to understand the technology before employing it in public sector contexts.
- **Governments should strongly consider establishing institutions, or empower existing institutions, to continuously oversee and publicly debate and define mandatory principles** to ensure that technology is developed, deployed, and used in the public interest.
- **Governments should ensure public funding of research on data practices and consumer and societal harms** stemming from generative AI.
- **International trade agreements must not hollow out transparency obligations** for generative AI systems, nor other obligations that are necessary to ensure consumer rights.
- **Shareholders and investors** in companies developing and deploying generative AI systems, particularly shareholder or investors from the public sector, should **require that steps are taken to avoid and/or mitigate exploitative practices, environmental impact, etc.** Companies should be required to have ethical guidelines as well as reporting on the steps taken.

#### 4.2.3 NEW LEGISLATIVE MEASURES

While there are already many legal frameworks in place that may be suited to address the harms of generative AI, there will undoubtedly be areas with legal gaps and loopholes. In cases where existing laws are not sufficient, it is necessary to create new frameworks to protect consumers from harm. As described in the previous chapter, there are already several legislative initiatives in progress, and it is crucial that these processes result in strong future proof rules that are founded on consumer and human rights.

We call on policymakers and lawmakers to take a strong stance in favour of consumer protection and preserving human rights. It is necessary to have robust legal measures, including strict obligations on developers and deployers of generative AI systems to operate in a transparent and accountable manner, and to restrict development, deployment, and use of systems that are fundamentally incompatible with these rights.

#### 4.2.3.1 Particular forms of generative AI that warrants additional scrutiny

- **Certain forms of manipulative techniques in generative AI systems should be banned.** This may include, for example, significant restrictions on anthropomorphized models, including the use of first-person language, the use of emojis and similar symbols, and simulating human emotions and similar attributes. Such restrictions could be dependent on the context and purpose of use. The threshold for acceptable techniques and applications should be higher when used by vulnerable groups such as children.
- For certain uses of generative AI systems, it might be necessary to **require pre-approval** from relevant enforcement agencies before deployment. Models that may lead to the exploitation or discrimination of consumers, in particular vulnerable consumers such as children, may be an example of this.
- **Policymakers must ensure that coming legislation is future proof**, to keep authorities from lagging behind rapid technological advancement. This involves principles and regulations that are technology neutral.

#### 4.2.3.2 Obligations for developers and deployers of generative AI

Responsible development, deployment and use of generative AI presupposes that it is possible to control how the systems work, inspect the training data, oversee social and environmental impacts, and more. While transparency itself is not a panacea, it is a prerequisite to ensure that technologies do not undermine consumer and human rights. This cannot be a responsibility left to the companies themselves.

There is an urgent need for independent oversight, research, and auditing of generative AI systems, to ensure that companies are held accountable if something goes wrong, to identify and root out biases and inaccuracies, and to otherwise ensure legal compliance and mitigate harms. Therefore, we present a number of measures that should be imposed on developers and deployers of generative AI systems.

### TRANSPARENCY

- Developers and deployers of generative AI systems must **be obligated to report and publish documentation about their risk assessments**, risk mitigation strategies, how they conduct content moderation, standardized performance metrics, etc. to the public. This should be done on two levels: a shorter, less technical version for consumers generally, as well as an in-depth description for experts from civil society, academia, and other third parties.
- All companies developing and deploying generative AI systems should have an **obligation to publish all information on energy use, water use and carbon emissions** for the whole lifecycle of the generative AI model and provide prognostics on future emissions for the day-to-day use. This includes resources needed for the production of the hardware, training of models, development, deployment, and use. A standardised model on calculating emissions, water use and energy use should be established that all companies need to use rather than inventing their own calculation system.
- Developers and deployers should disclose the names of all their suppliers publicly and report in **a transparent way on the working conditions in their whole supply chain**, including living wage and psychological support for moderators of violent and disturbing content, and packages for labourers needed only temporarily for a given task.

- Developers of baseline models for generative AI must be obligated to register their models in a centralized, public system, to ensure oversight of relevant generative AI models.
- Deployers of generative AI in consumer facing interfaces and services **must be obligated to disclose how the generated content is influenced by commercial interests of developers, deployers, or third parties**. This is particularly relevant when the content generated serves to inform consumer choices, such as content generated in the context of search queries or similar.
- Deployers of generative AI systems must be **obligated to disclose whenever consumers are interacting with a generative AI system**, and whether consumer-facing systems are using artificial intelligence to affect the outcome of a decision.
- Public and private organizations must be obligated to **disclose whenever content has been generated by generative AI**, when that content may have an effect on decisions affecting consumers, consumer rights more broadly, or democratic processes.

#### RISK MITIGATION

- Deployers of generative AI systems must be obligated to carefully consider the context in which the generative AI system is to be deployed. Deployers of generative AI systems **should not use generative AI systems without a careful risk assessment**, including a mapping of what problems the system is meant to solve, verification that the system is compliant with relevant laws, risks to consumers and consumer rights, risks to human rights, risks of harm likely to impact vulnerable groups, adverse impacts on the environment, foreseeable societal and collective harms, privacy harms, etc.
- Deployers of generative AI systems must be **obligated to implement effective measures to mitigate the risks uncovered** in the risk assessment before deploying the system, to arrive at acceptable residual risk. If the risks cannot be mitigated, or the system does not solve the problems it is meant to solve, the system should not be deployed in that context.
- Developers and deployers of generative AI systems must be **obligated to involve representatives from populations that may be impacted by the technology**, in particular marginalized and vulnerable groups and communities. This precipitates democratic participation and interdisciplinary involvement. Stakeholder participation is necessary in the context of development and training of generative AI models, and for associated topics such as risk assessments, risk mitigation strategies, and content moderation, which needs to account for different cultural contexts, languages, etc.
- Deployers of generative AI systems must be **obligated to monitor and address the system's impact on consumers** after deploying the system, conducting continuous risk assessments and mitigation to arrive at acceptable residual risk, taking particular account of impacts on marginalized and vulnerable groups and communities.

## ACCOUNTABILITY

- There must be **clear rules on accountability and liability for harmful effects of generative AI systems**, such as harms to privacy, safety, consumer rights, and fundamental rights more broadly. These rules must clearly indicate which company in the supply chain is liable or require that developers and deployers of generative AI clearly establish responsibility and liability between themselves.
- Any accountability scheme must **make it easy for consumers, enforcement agencies and courts to hold the companies liable for consumer harm**.
- Developers of generative AI systems must be **responsible for the data they use**, representativeness in the data sets, their data cleaning and labelling practices, and other design choices that will affect all downstream uses of the systems. Such choices must be carefully documented, so downstream developers and deployers can consider the generative AI system's risks and suitability.
- **Technical standards and certification schemes should be developed** and used to assist developers and deployers of generative AI systems to develop, train, deploy and use the system in a responsible and legal way. However, policymakers must not outsource human rights, political and legal issues to standards bodies. Governments must ensure civil society participation in such bodies, or, lacking such participation, should not rely on them.
- Generative AI systems and models must be **auditable by independent researchers, enforcement agencies, and other third parties**. This is essential to mitigate the risk of bias and discrimination, ensure responsible use of training data, and to ensure compliance with applicable legal requirements.
- **Audits should at least involve the training data, data collection practices, data labelling practices, content moderation practices, sustainability reports and the algorithmic models**. Audits should be carefully documented to ensure accountability and reproducibility, and should be based on standardized auditing requirements.
- Companies should be **obligated to have quantitative and timebound commitments to reduce consumption**, based on the calculations of the carbon emissions, energy, and water use of developing and deploying generative AI. This progress should also be audited by an external and independent actor with public reporting. This means bigger models might be reduced and less ambitious. Claims of zero net carbon activities, and carbon offsetting schemes should not be the default model with which companies "compensate" emissions, they should reduce emissions in their own activities.



## Endnotes

- 1 “ChatGPT reaches 100 million users two months after launch”, Dan Milmo, The Guardian (2023). <https://www.theguardian.com/technology/2023/feb/02/chatgpt-100-million-users-open-ai-fastest-growing-app>
- 2 Large language models are sophisticated AI models that are designed to generate text that resembles human language. They are normally trained on vast amounts of text sources to “learn” patterns and grammar. LLMs can be used for tasks such as machine translation, sentiment analysis, human-machine interaction, proofreading, and many other purposes.
- 3 ChatGPT, <https://chat.openai.com/>
- 4 “Microsoft reportedly to add ChatGPT to Bing search engine”, Johana Bhuiyan, The Guardian (2023). <https://www.theguardian.com/technology/2023/jan/05/microsoft-chatgpt-bing-search-engine>
- 5 “Microsoft rolls out ChatGPT-powered Teams Premium”, Reuters (2023). <https://www.reuters.com/technology/microsoft-rolls-out-chatgpt-powered-teams-premium-2023-02-02/>
- 6 “A New Chat Bot Is a ‘Code Red’ for Google’s Search Business”, The New York Times, (2023). <https://www.nytimes.com/2022/12/21/technology/ai-chatgpt-google-search.html>
- 7 “A new era for AI and Google Workspace”, Product Announcement, Google (2023). <https://workspace.google.com/blog/product-announcements/generative-ai>
- 8 “The Galactica AI model was trained on scientific knowledge – but it spat out alarmingly plausible nonsense”, Aaron J. Snowwell, Jean Burgess (2022). <https://theconversation.com/the-galactica-ai-model-was-trained-on-scientific-knowledge-but-it-spat-out-alarmingly-plausible-nonsense-195445>
- 9 “Facebook’s Powerful Large Language Model Leaks Online”, Joseph Cox, The Vice (2023). <https://www.vice.com/en/article/xgwqgw/facebook-powerful-large-language-model-leaks-online-4chan-llama>
- 10 Hugging Face. <https://huggingface.co/bigscience/bloom>
- 11 “Stability AI Launches the First of its StableLM Suite of Language Models”, Stability AI (2023). <https://stability.ai/blog/stability-ai-launches-the-first-of-its-stablelm-suite-of-language-models>
- 12 Midjourney. <https://www.midjourney.com/>
- 13 Dall-E. <https://labs.openai.com>
- 14 Stability AI. <https://stability.ai>
- 15 MusicStar.ai. <https://musicstar.ai>
- 16 “Herzog and Žižek become uncanny AI bots trapped in endless conversation”, Benj Edwards, Ars Technica (2022). <https://arstechnica.com/information-technology/2022/11/herzog-and-zizek-become-uncanny-ai-bots-trapped-in-endless-conversation/>
- 17 ElevenLabs. <https://beta.elevenlabs.io>
- 18 Vall-E. <https://vall-e.io>
- 19 “Introducing Make-A-Video: An AI system that generates videos from text”. Meta, (2022). <https://ai.facebook.com/blog/generative-ai-text-to-video/>
- 20 Make-A-Video. <https://makeavideo.studio>
- 21 Dreamix. <https://dreamix-video-editing.github.io>
- 22 Stability Animation. <https://platform.stability.ai/docs/features/animation>
- 23 “Create generative AI video-to-video right from your phone with Runway’s iOS app”, James Vincent, The Verge, (2023). <https://www.theverge.com/2023/4/24/23695788/generative-ai-video-runway-mobile-app-ios>
- 24 “Call for action to open an inquiry on generative AI systems to address risks and harms for consumers”, BEUC (2023). [https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023-045\\_Call\\_for\\_action\\_CPC\\_authorities\\_Generative\\_AI\\_systems.pdf](https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023-045_Call_for_action_CPC_authorities_Generative_AI_systems.pdf)
- 25 “End of the Billable Hour? Law Firms Get On Board With Artificial Intelligence”, Erin Mulvaney, Lauren Weber, The Washington Post (2023). <https://www.wsj.com/articles/end-of-the-billable-hour-law-firms-get-on-board-with-artificial-intelligence-17ebd3f8>
- 26 “ELIZA—A Computer Program For the Study of Natural Language Communication Between Man and Machine”, Joseph Weizenbaum (1966). [http://www.universelle-automation.de/1966\\_Boston.pdf](http://www.universelle-automation.de/1966_Boston.pdf)
- 27 “Pandora’s Box: Generative AI Companies, ChatGPT, and Human Rights”, Human Rights Watch (2023). <https://www.hrw.org/news/2023/05/03/pandoras-box-generative-ai-companies-chatgpt-and-human-rights>
- 28 “AI machines aren’t ‘hallucinating’. But their makers are”, Naomi Klein, The Guardian (2023). <https://www.theguardian.com/commentisfree/2023/may/08/ai-machines-hallucinating-naomi-klein>
- 29 “How to Recognize AI Snake Oil”, Doug Hulet (2021). <https://www.cs.princeton.edu/news/how-recognize-ai-snake-oil>
- 30 “Pause Giant AI Experiments: An Open Letter”, Future of life Institute (2023). <https://futureoflife.org/open-letter/pause-giant-ai-experiments/>
- 31 “Policy makers: Please don’t fall for the distractions of #Alhype”, Emily M. Bender, (2023). <https://medium.com/@emilymenon-bender/policy-makers-please-dont-fall-for-the-distractions-of-aihype-e03fa80ddb1>
- 32 “The so-called ‘Godfather of the A.I.’ joins The Lead to offer a dire warning about the dangers of artificial intelligence”, Geoffrey Hinton, CNN (2023). <https://edition.cnn.com/videos/tv/2023/05/02/the-lead-geoffrey-hinton.cnn>
- 33 “What you need to know about generative AI and human rights”, Access Now (2023). <https://www.accessnow.org/what-you-need-to-know-about-generative-ai-and-human-rights/>
- 34 “The folly of technological solutionism: An interview with Evgeny Morozov”. Natasha Dow Schüll, Public Books (2013). <https://www.publicbooks.org/the-folly-of-technological-solutionism-an-interview-with-evgeny-morozov/>
- 35 “The Delusion at the Center of the A.I. Boom”, Evan Selinger, The Slate, (2023). <https://slate.com/technology/2023/03/>

- [chatgpt-artificial-intelligence-solution-ism-hype.html](#)
- 36 “Generative AI and the Digital Commons”, Saffron Huang, Divya Siddarth, (2023). <https://arxiv.org/abs/2303.11074>
- 37 “Indigenous groups in NZ, US fear colonisation as AI learns their languages”, Rina Chandran, Context (2023). <https://www.context.news/ai/nz-us-indigenous-fear-colonisation-as-bots-learn-their-languages>
- 38 “Lina Khan: We Must Regulate A.I. Here’s How.”, Lina Khan, New York Times, (2023). <https://www.nytimes.com/2023/05/03/opinion/ai-lina-khan-ftc-technology.html>
- 39 “What does a leaked Google memo reveal about the future of AI?”, The Economist (2023). <https://www.economist.com/leaders/2023/05/11/what-does-a-leaked-google-memo-reveal-about-the-future-of-ai>
- 40 “The Imminent Danger of A.I. Is One We’re Not Talking About”, Ezra Klein, The New York Times (2023). <https://www.nytimes.com/2023/02/26/opinion/microsoft-bing-sydney-artificial-intelligence.html>
- 41 “Microsoft pushing you to set Bing and Edge as your defaults to get its new OpenAI-powered search engine faster is giving off big 1990s energy”, Hasan Chowdhury and Shona Ghosh, Insider (2023). <https://www.businessinsider.com/microsoft-wants-to-repeat-1990s-dominance-with-new-bing-ai-2023-2>
- 42 “Walled garden”, Andrew Froehlich, TechTarget (2023). <https://www.techtarget.com/searchsecurity/definition/walled-garden>
- 43 “Snapchat’s AI chatbot is now free for all global users, says the AI will later ‘Snap’ you back”, Sarah Perez, TechCrunch (2023). <https://techcrunch.com/2023/04/19/snapchat-opens-its-ai-chatbot-to-global-users-says-the-ai-will-later-snap-you-back/>
- 44 “Publishers Worry A.I. Chatbots Will Cut Readership”. Katie Robertson, The New York Times (2023). <https://www.nytimes.com/2023/03/30/business/media/publishers-chatbots-search-engines.html>
- 45 “The AI takeover of Google Search starts now”, David Pierce, The Verge (2023). <https://www.theverge.com/2023/5/10/23717120/google-search-ai-results-generated-experience-io>
- 46 “Indigenous groups in NZ, US fear colonisation as AI learns their languages”, Rina Chandran, Context (2023). <https://www.context.news/ai/nz-us-indigenous-fear-colonisation-as-bots-learn-their-languages>
- 47 “Google shared AI knowledge with the world – until ChatGPT caught up”, Nitasha Tiku and Gerrit De Vynck, The Washington Post (2023). <https://www.washingtonpost.com/technology/2023/05/04/google-ai-stop-sharing-research/>
- 48 Hypothetical AI system that demonstrates human-level intelligence and autonomy. Does not currently exist.
- 49 “Microsoft Says New A.I. Shows Signs of Human Reasoning”, Cade Metz, The New York Times (2023). <https://www.nytimes.com/2023/05/16/technology/microsoft-ai-human-reasoning.html>
- 50 “OpenAI co-founder on company’s past approach to openly sharing research: ‘We were wrong’”, James Vincent, The Verge (2023). <https://www.theverge.com/2023/3/15/23640180/openai-gpt-4-launch-closed-research-ilya-sutskever-interview>
- 51 “We are a little bit scared’: OpenAI CEO warns of risks of artificial intelligence”, Edward Helmore, The Guardian (2023). <https://www.theguardian.com/technology/2023/mar/17/openai-sam-altman-artificial-intelligence-warning-gpt4>
- 52 “GPT-4 and professional benchmarks: the wrong answer to the wrong question”, Arvind Narayanan and Sayash Kapoor, AI Snake Oil (2023). <https://aisnakeoil.substack.com/p/gpt-4-and-professional-benchmarks>
- 53 “OpenAI’s policies hinder reproducible research on language models”, Arvind Narayanan and Sayash Kapoor, AI Snake Oil (2023). <https://aisnakeoil.substack.com/p/openais-policies-hinder-reproducible-research-on-language-models>
- 54 “How trade commitments narrowed EU rules to access AI’s source codes”, Luca Bertuzzi, Euractiv (2023). <https://www.euractiv.com/section/artificial-intelligence/news/how-trade-commitments-narrowed-eu-rules-to-access-ais-source-codes/>
- 55 “Early thoughts on regulating generative AI like ChatGPT”, Alex Engler, Brookings. (2023). <https://www.brookings.edu/blog/techtank/2023/02/21/early-thoughts-on-regulating-generative-ai-like-chatgpt/>
- 56 “Klarna brings smooth shopping to ChatGPT”, Klarna (2023). <https://www.klarna.com/international/press/klarna-brings-smooth-shopping-to-chatgpt/>
- 57 “EU Consumer Protection 2.0 - Protecting fairness and consumer choice in a digital economy”, BEUC (2022). [https://www.beuc.eu/sites/default/files/publications/beuc-x-2022-015\\_protecting\\_fairness\\_and\\_consumer\\_choice\\_in\\_a\\_digital\\_economy.pdf](https://www.beuc.eu/sites/default/files/publications/beuc-x-2022-015_protecting_fairness_and_consumer_choice_in_a_digital_economy.pdf)
- 58 “What Really Happened When Google Ousted Timnit Gebru”, Tom Simonite, Wired, (2022). <https://www.wired.com/story/google-timnit-gebru-ai-what-really-happened/>
- 59 “Big tech companies cut AI ethics staff, raising safety concerns”, Cristina Criddle and Madhumita Murgia, Financial Times (2023). <https://www.ft.com/content/26372287-6fb3-457b-9e9c-f722027f36b3>
- 60 “ChatGPT maker OpenAI calls for AI regulation, warning of ‘existential risk’”, Ellen Francis, The Washington Post, (2023). <https://www.washingtonpost.com/technology/2023/05/24/chatgpt-openai-artificial-intelligence-regulation/>
- 61 “OpenAI may leave the EU if regulations bite”, Reuters (2023). <https://www.reuters.com/technology/openai-may-leave-eu-if-regulations-bite-ceo-2023-05-24/>
- 62 “Help! My Political Beliefs Were Altered by a Chatbot!”, Christopher Mims, The Wall Street Journal (2023). <https://www.wsj.com/articles/chatgpt-bard-bing-ai-political-beliefs-151a0fe4>
- 63 “Stack Overflow Bans ChatGPT For Constantly Giving Wrong Answers”, Janus Rose, Vice (2022). <https://www.vice.com/en/article/wxnaem/stack-overflow-bans-chatgpt-for-constantly-giving-wrong-answers>
- 64 “AI-assisted plagiarism? ChatGPT bot says it has an answer for that”, Alex Hern, The Guardian (2022). <https://www.theguardian.com/technology/2022/dec/31/ai-assisted-plagiarism-chatgpt-bot-says-it-has-an-an->

- [swer-for-that](#)
- 65 “Google employees label AI chatbot Bard ‘worse than useless’ and ‘a pathological liar’”, James Vincent, *The Verge* (2023). <https://www.theverge.com/2023/4/19/23689554/google-ai-chatbot-bard-employees-criticism-pathological-liar>
- 66 “ChatGPT is making up fake Guardian articles. Here’s how we’re responding”, Chris Moran, *The Guardian* (2023). <https://www.theguardian.com/commentisfree/2023/apr/06/ai-chatgpt-guardian-technology-risks-fake-article>
- 67 “Publishers tout generative AI opportunities to save and make money amid rough media market”, Sara Guaglione, *Digday* (2023). <https://digiday.com/media/publishers-tout-generative-ai-opportunities-to-save-and-make-money-amid-rough-media-market/>
- 68 “CNET Is Reviewing the Accuracy of All Its AI-Written Articles After Multiple Major Corrections”, Lauren Leffer (2023). <https://gizmodo.com/cnet-ai-chatgpt-news-robot-1849996151>
- 69 “Situating Search”. Chirag Shah and Emily Bender (2022). <https://dl.acm.org/doi/pdf/10.1145/3498366.3505816>
- 70 “People Are Using AI for Therapy, Even Though ChatGPT Wasn’t Built for It”, Rachel Metz, *Bloomberg* (2023). <https://www.bloomberg.com/news/articles/2023-04-18/ai-therapy-becomes-new-use-case-for-chatgpt>
- 71 “Governments are embracing ChatGPT-like bots. Is it too soon?”, J.D. Capelouto and Diego Mendoza, *Semafor* (2023). <https://www.semafor.com/article/03/03/2023/governments-using-chatgpt-bots>
- 72 “The dark side of artificial intelligence: manipulation of human behaviour”, Georgios Petropoulos, *Bruegel* (2023). <https://www.bruegel.org/blog-post/dark-side-artificial-intelligence-manipulation-human-behaviour>
- 73 “People keep anthropomorphizing AI. Here’s why”, Arvind Narayanan and Sayash Kapoor, *AI Snake Oil* (2023). <https://aisnakeoil.substack.com/p/people-keep-anthropomorphizing-ai>
- 74 “We warned Google that people might believe AI was sentient. Now it’s happening”, Timnit Gebru, Margaret Mitchell, *Washington Post* (2022). <https://www.washingtonpost.com/opinions/2022/06/17/google-ai-ethics-sentient-lemoine-warning/>
- 75 “Column: Afraid of AI? The startups selling it want you to be”, Brian Merchant *Los Angeles Times* (2023). <https://www.latimes.com/business/technology/story/2023-03-31/column-afraid-of-ai-the-startups-selling-it-want-you-to-be>
- 76 “AI Doesn’t Hallucinate. It Makes Things Up”, Rachel Metz, *Bloomberg* (2023). <https://www.bloomberg.com/news/newsletters/2023-04-03/chatgpt-bing-and-bard-don-t-hallucinate-they-fabricate>
- 77 “Chatbots shouldn’t use emojis”, Carissa Véliz (2023). <https://www.nature.com/articles/d41586-023-00758-y>
- 78 “The engineer who claimed a Google AI is sentient has been fired”, Mitchell Clark, *The Verge* (2022). <https://www.theverge.com/2022/7/22/23274958/google-ai-engineer-blake-lemoine-chatbot-lambda-2-sentience>
- 79 “Users Report Microsoft’s ‘Unhinged’ Bing AI Is Lying, Berating Them”, Jordan Pearson, *Vice* (2023). <https://www.vice.com/en/article/3ad39b/microsoft-bing-ai-unhinged-lying-berating-users>
- 80 “On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?”, Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, Shmargaret Shmitchell [sic] (2023). <https://dl.acm.org/doi/pdf/10.1145/3442188.3445922>
- 81 “Human heuristics for AI-generated language are flawed”, Maurice Jakesch, Jeffrey T. Hancock and Mor Naaman (2023). <https://www.pnas.org/doi/10.1073/pnas.2208839120>
- 82 “Privacy Concerns in Chatbot Interactions”, Carolin Ischen, Theo Araujo, Hilde Voorveld, Guda van Noord and Edith Smit (2020). [https://link.springer.com/chapter/10.1007/978-3-030-39540-7\\_3](https://link.springer.com/chapter/10.1007/978-3-030-39540-7_3)
- 83 “It’s Hurting Like Hell’: AI Companion Users Are In Crisis, Reporting Sudden Sexual Rejection”, Samantha Cole, *Vice* (2023). <https://www.vice.com/en/article/y3py9j/ai-companion-replika-erotic-roleplay-updates>
- 84 “Snapchat is releasing its own AI chatbot powered by ChatGPT”, Alex Heath, *The Verge* (2023). <https://www.theverge.com/2023/2/27/23614959/snapchat-my-ai-chatbot-chatgpt-openai-plus-subscription>
- 85 “Snapchat tried to make a safe AI. It chats with me about booze and sex”, Geoffrey A. Fowler, *The Washington Post* (2023). <https://www.washingtonpost.com/technology/2023/03/14/snapchat-myai/>
- 86 “Emotional artificial intelligence in children’s toys and devices: Ethics, governance and practical remedies”, Andrew McStay and Gilad Rosne (2021). <https://journals.sagepub.com/doi/10.1177/2053951721994877>
- 87 “Deepfakes for all: Uncensored AI art model prompts ethics questions”, Kyle Wiggers, *TechCrunch* (2022). <https://techcrunch.com/2022/08/24/deepfakes-for-all-uncensored-ai-art-model-prompts-ethics-questions/>
- 88 “Facing reality? Law enforcement and the challenge of deepfakes”, *Europol Innovation Lab* (2022). [https://www.europol.europa.eu/cms/sites/default/files/documents/Europol\\_Innovation\\_Lab\\_Facing\\_Reality\\_Law\\_Enforcement\\_And\\_The\\_Challenge\\_Of\\_Deepfakes.pdf](https://www.europol.europa.eu/cms/sites/default/files/documents/Europol_Innovation_Lab_Facing_Reality_Law_Enforcement_And_The_Challenge_Of_Deepfakes.pdf)
- 89 “Fake images of Trump arrest show ‘giant step’ for AI’s disruptive power”, Isaac Stanley-Becker, Naomi Nix, *The Washington Post* (2023). <https://www.washingtonpost.com/politics/2023/03/22/trump-arrest-deep-fakes/>
- 90 “Found through Google, bought with Visa and Mastercard: Inside the deepfake porn economy”, Kat Tenbarge, *NBC* (2023). <https://www.nbcnews.com/tech/internet/deepfake-porn-ai-mr-deep-fake-economy-google-visa-mastercard-download-rcna75071>
- 91 “Scammers are using AI to impersonate your loved ones. Here’s what to watch out for”, Sabrina Ortiz, *CDNET* (2023).
- 92 “Exclusive: GPT-4 readily spouts misinformation, study finds”, Sara Fischer, *Axios* (2023). <https://www.axios.com/2023/03/21/gpt4-misinformation-newsguard-study>
- 93 “AI presents political peril for 2024

- with threat to mislead voters”, David Klepper and Ali Swenson, AP News (2023). <https://apnews.com/article/artificial-intelligence-misinformation-deepfakes-2024-election-trump-59fb51002661ac-5290089060b3ae39a0>
- 94 “Watermarking ChatGPT, DALL-E and other generative AIs could help protect against fraud and misinformation”, Hany Farid. *The Conversation* (2023). <https://theconversation.com/watermarking-chatgpt-dall-e-and-other-generative-ais-could-help-protect-against-fraud-and-misinformation-202293>
- 95 “Google introduces new features to help identify AI images in Search and elsewhere”, Sarah Perez, TechCrunch (2023). <https://techcrunch.com/2023/05/10/google-introduces-new-features-to-help-identify-ai-images-in-search-and-elsewhere/>
- 96 “Lærere fortvilet over ny kunstig intelligens”, Daniel Eriksen, NRK (2022). <https://www.nrk.no/kultur/laerere-fortvilet-over-ny-kunstig-intelligens-1.16210580>
- 97 “OpenAI’s attempts to watermark AI text hit limits”, Kyle Wiggers, TechCrunch (2022). <https://techcrunch.com/2022/12/10/openais-attempts-to-watermark-ai-text-hit-limits/>
- 98 “We pitted ChatGPT against tools for detecting AI-written text, and the results are troubling”, Armin Alimardani and Emma A. Jane, *The Conversation* (2023). <https://theconversation.com/we-pitted-chatgpt-against-tools-for-detecting-ai-written-text-and-the-results-are-troubling-199774>
- 99 “OpenAI releases tool to detect AI-generated text, including from ChatGPT”, Kyle Wiggers, TechCrunch (2023). <https://techcrunch.com/2023/01/31/openai-releases-tool-to-detect-ai-generated-text-including-from-chatgpt/>
- 100 “ChatGPT and Generative AI in Content Marketing”, Kelsey Voss, *Insider Intelligence* (2023). <https://www.insiderintelligence.com/content/chatgpt-generative-ai-content-marketing>
- 101 “The Generative AI Revolution Is Creating The Next Phase Of Autonomous Enterprise”, Mark Minevich, *Forbes* (2023). <https://www.forbes.com/sites/markminevich/2023/01/29/the-generative-ai-revolution-is-creating-the-next-phase-of-autonomous-enterprise/>
- 102 “The Generative AI Revolution Is Creating The Next Phase Of Autonomous Enterprise”, Mark Minevich, *Forbes* (2023). <https://adage.com/article/marketing-news-strategy/levis-uses-ai-models-increase-diversity-in-cites-backlash/2482046>
- 103 “Beck’s uses AI to create, and advertise, a new beer in experiment for future use cases”, Hannah Bowler, *The Forbes* (2023). <https://www.thedrum.com/news/2023/03/30/beck-s-uses-ai-create-and-advertise-new-beer-experiment-future-use-cases>
- 104 “Ads are coming for the Bing AI chatbot, as they come for all Microsoft products”, Andrew Cunningham, *ArsTechnica* (2023). <https://arstechnica.com/gadgets/2023/03/ads-are-coming-for-the-bing-ai-chatbot-as-they-come-for-all-microsoft-products/>
- 105 “Yes, Google’s AI-infused search engine will have ads”, Ryan Barwick, *Marketing Brew* (2023). <https://www.marketingbrew.com/stories/2023/05/23/yes-google-s-ai-infused-search-engine-will-have-ads>
- 106 “The dark side of artificial intelligence: manipulation of human behaviour”, Georgios Petropoulos, *Bruegel* (2023). <https://www.bruegel.org/blog-post/dark-side-artificial-intelligence-manipulation-human-behaviour>
- 107 The Norwegian Consumer Council has previously described the negative effects of surveillance-based advertising and has called for a general ban (2021). <https://www.forbrukerradet.no/side/new-report-details-threats-to-consumers-from-surveillance-based-advertising/>
- 108 See for example Shoshana Zuboff’s book “The Age of Surveillance Capitalism – The Fight for a Human Future at the New Frontier of Power” (2019).
- 109 “Exploring 12 Million of the 2.3 Billion Images Used to Train Stable Diffusion’s Image Generator”, Ando Baio, *Waxy* (2022). <https://waxy.org/2022/08/exploring-12-million-of-the-images-used-to-train-stable-diffusion-image-generator/>
- 110 “Artist finds private medical record photos in popular AI training data set”, Benj Edwards (*ArsTechnica*), 2022. <https://arstechnica.com/information-technology/2022/09/artist-finds-private-medical-record-photos-in-popular-ai-training-data-set/>
- 111 “On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?”, Emily M. Bender, Timnit Gebru, Angelina McMillan-Major and Shmargaret Shmitchell [sic] (2021). <https://dl.acm.org/doi/pdf/10.1145/3442188.3445922>
- 112 “On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?”, Emily M. Bender, Timnit Gebru, Angelina McMillan-Major and Shmargaret Shmitchell [sic] (2021). <https://dl.acm.org/doi/pdf/10.1145/3442188.3445922>
- 113 “The viral AI avatar app Lensa undressed me—without my consent”, Melissa Heikkilä, *MIT Technology Review* (2022). <https://www.technologyreview.com/2022/12/12/1064751/the-viral-ai-avatar-app-lensa-undressed-me-without-my-consent/>
- 114 “Stable Diffusion and DALL-E display bias when prompted for artwork of ‘African workers’ versus ‘European workers’”, Thomas Maxwell, *Insider* (2023). <https://www.businessinsider.com/ai-image-prompt-for-african-workers-depicts-harmful-stereotypes-2023-4>
- 115 “Inside the secret list of websites that make AI like ChatGPT sound smart”, Kevin Schaul, Szu Yu Chen and Nitasha Tiku, *The Washington Post* (2023). <https://www.washingtonpost.com/technology/interactive/2023/ai-chatbot-learning/>
- 116 “Google ‘fixed’ its racist algorithm by removing gorillas from its image-labeling tech”, James Vincent, *The Verge* (2018). <https://www.theverge.com/2018/1/12/16882408/google-racist-gorillas-photo-recognition-algorithm-ai>
- 117 “Facebook apology as AI labels black men primates”, *BBC* (2021). <https://www.bbc.com/news/technology-58462511>
- 118 “On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?”,

- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major and Shmargaret Shmitchell [sic] (2021). <https://dl.acm.org/doi/pdf/10.1145/3442188.3445922>
- 119 “There’s More to AI Bias Than Biased Data, NIST Report Highlights”, The National Institute of Standards and Technology (2022). <https://www.nist.gov/news-events/news/2022/03/theres-more-ai-bias-biased-data-nist-report-highlights>
- 120 “Ensuring artificial intelligence (AI) technologies for health benefit older people”, Dr Vânia de la Fuente-Núñez, WHO (2022). [https://www.who.int/news/item/09-02-2022-ensuring-artificial-intelligence-\(ai\)-technologies-for-health-benefit-older-people](https://www.who.int/news/item/09-02-2022-ensuring-artificial-intelligence-(ai)-technologies-for-health-benefit-older-people)
- 121 “AI image generator Midjourney blocks porn by banning words about the human reproductive system”, Melissa Heikkilä, MIT Technology Review (2023). <https://www.technologyreview.com/2023/02/24/1069093/ai-image-generator-midjourney-blocks-porn-by-banning-words-about-the-human-reproductive-system/>
- 122 “How a tiny company with few rules is making fake images go mainstream”, Isaac Stanley-Becker, The Washington Post (2023). <https://www.washingtonpost.com/technology/2023/03/30/midjourney-ai-image-generation-rules/>
- 123 “Oh No, ChatGPT AI Has Been Jailbroken To Be More Reckless”, Claire Jackson, Kotaku (2023). <https://kotaku.com/chatgpt-ai-openai-dan-censorship-chatbot-red-dit-1850088408>
- 124 “Quantifying ChatGPT’s gender bias”, Sayash Kapoor and Arvind Narayanan. AI Snake Oil (2023). <https://aisnakeoil.substack.com/p/quantifying-chatgpts-gender-bias>
- 125 “AI Is Steeped in Big Tech’s Digital Colonialism”, Grace Browne, Wired (2023). <https://www.wired.com/story/abeba-birhane-ai-datasets/>
- 126 “The AI chatbot can misrepresent key facts with great flourish, even citing a fake Washington Post article as evidence”, Pranshu Verma and Will Oremus, The Washington Post (2023). <https://www.washingtonpost.com/technology/2023/04/05/chatgpt-lies/>
- 127 “How I Broke Into a Bank Account With an AI-Generated Voice”, Joseph Cox, Vice, (2023). <https://www.vice.com/en/article/dy7axa/how-i-broke-into-a-bank-account-with-an-ai-generated-voice>
- 128 “Scammers use AI to enhance their family emergency schemes”, FTC (2023). <https://consumer.ftc.gov/consumer-alerts/2023/03/scammers-use-ai-enhance-their-family-emergency-schemes>
- 129 “Three ways AI chatbots are a security disaster”, Melissa Heikkilä, MIT Technology Review (2023). <https://www.technologyreview.com/2023/04/03/1070893/three-ways-ai-chatbots-are-a-security-disaster/>
- 130 “Opwnai: Cybercriminals starting to use chatgpt”, Check Point Research (2023). <https://research.checkpoint.com/2023/opwnai-cybercriminals-starting-to-use-chatgpt/>
- 131 “Dual use of artificial-intelligence-powered drug discovery”, Fabio Urbina, Filippa Lentzos, Cédric Invernizzi and Sean Ekins (2022). <https://www.nature.com/articles/s42256-022-00465-9.epdf>
- 132 “ChatGPT - The impact of Large Language Models on Law Enforcement”, Europol (2023). <https://www.europol.europa.eu/publications-events/publications/chatgpt-impact-of-large-language-models-law-enforcement>
- 133 “Companies are struggling to keep corporate secrets out of ChatGPT”, Sam Sabin, Axios (2023). <https://www.axios.com/2023/03/10/chatgpt-ai-cybersecurity-secrets>
- 134 “Amazon warns employees not to share confidential information with ChatGPT after seeing cases where its answer ‘closely matches existing material’ from inside the company”, Eugene Kim, Business Insider (2023). <https://www.businessinsider.com/amazon-chatgpt-openai-warns-employees-not-share-confidential-information-microsoft-2023-1>
- 135 “OpenAI ceo says ai will give medical advice to people too poor to afford doctors”, Frank Landmore, The Byte (2021). <https://futurism.com/the-byte/openai-ceo-ai-medical-advice>
- 136 “Helpline workers for the National Eating Disorder Association say they are being replaced by AI”, Britney Nguyen, Insider (2023). <https://www.businessinsider.com/eating-disorders-nonprofit-reportedly-fired-humans-offer-ai-chatbot-2023-5>
- 137 See for example the Commission’s draft of the AIA art. 14(1), [https://eur-lex.europa.eu/resource.html?uri=cellar:e0649735-a372-11eb-9585-01aa75ed71a1.0001.02/DOC\\_1&format=PDF](https://eur-lex.europa.eu/resource.html?uri=cellar:e0649735-a372-11eb-9585-01aa75ed71a1.0001.02/DOC_1&format=PDF).
- 138 See some examples of this in relation to art. 22 GDPR, where companies have effectively been rubberstamping decisions. “Automated Decision-Making Under the GDPR: Practical Cases from Courts and Data Protection Authorities”, Future of Privacy Forum (2022). <https://fpf.org/wp-content/uploads/2022/05/FPF-ADM-Report-R2-singles.pdf>.
- 139 Also called “automation bias” (see for example “Automation Bias in Intelligent Time Critical Decision Support Systems”, M.L. Cumming (2012), <https://arc.aiaa.org/doi/10.2514/6.2004-6313>) and “algorithm aversion” (see for example “Algorithm Aversion: People Erroneously Avoid Algorithms after Seeing Them Err” Dietvorst, Simmons and Massey (2014) [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2466040](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2466040)), respectively.
- 140 “How Can Artificial Intelligence Combat Climate Change?”, World101. <https://world101.cfr.org/global-era-issues/climate-change/how-can-artificial-intelligence-combat-climate-change> and “AI Is Essential for Solving the Climate Crisis”, Hamid Maher, Hubertus Meinecke, Damien Gromier, Mateo Garcia-Novelli and Ruth Fortmann, Boston Consulting Group (2022). <https://www.bcg.com/publications/2022/how-ai-can-help-climate-change>
- 141 “Emissions Gap Report 2022”, United Nations Environment Programme (2022). <https://www.unep.org/news-and-stories/story/new-pact-tech-companies-take-climate-change>
- 142 “Artificial Intelligence Is Booming—So Is Its Carbon Footprint”, Josh Saul and Dina Bass, Bloomberg (2023). <https://www.bloomberg.com/news/articles/2023-03-09/how-much->

- energy-do-ai-and-chatgpt-use-no-one-knows-for-sure
- 143 “Generative AI Breaks The Data Center: Data Center Infrastructure And Operating Costs Projected To Increase To Over \$76 Billion By 2028”, Jim McGregor, Forbes (2023). <https://www.forbes.com/sites/tiriasresearch/2023/05/12/generative-ai-breaks-the-data-center-data-center-infrastructure-and-operating-costs-projected-to-increase-to-over-76-billion-by-2028/>
- 144 “Generative AI Breaks The Data Center: Data Center Infrastructure And Operating Costs Projected To Increase To Over \$76 Billion By 2028”, Jim McGregor, Forbes (2023). <https://www.forbes.com/sites/tiriasresearch/2023/05/12/generative-ai-breaks-the-data-center-data-center-infrastructure-and-operating-costs-projected-to-increase-to-over-76-billion-by-2028/>
- 145 “Generative AI Breaks The Data Center: Data Center Infrastructure And Operating Costs Projected To Increase To Over \$76 Billion By 2028”, Jim McGregor, Forbes (2023). <https://www.forbes.com/sites/tiriasresearch/2023/05/12/generative-ai-breaks-the-data-center-data-center-infrastructure-and-operating-costs-projected-to-increase-to-over-76-billion-by-2028/>
- 146 “Generative AI Breaks The Data Center: Data Center Infrastructure And Operating Costs Projected To Increase To Over \$76 Billion By 2028”, Jim McGregor, Forbes (2023). <https://www.forbes.com/sites/tiriasresearch/2023/05/12/generative-ai-breaks-the-data-center-data-center-infrastructure-and-operating-costs-projected-to-increase-to-over-76-billion-by-2028/>
- 147 “The Generative AI Race Has a Dirty Secret”, Chris Stokel-Walker, Wired (2023). <https://www.wired.com/story/the-generative-ai-search-race-has-a-dirty-secret/>
- 148 “Generating Harms: Generative AI’s Impact & Paths Forward”, EPIC (2023). [EPIC-Generative-AI-White-Paper-May2023.pdf](https://www.epic-generative-ai-white-paper-may2023.pdf) (p. 40)
- 149 “AI and the Challenge of Sustainability”, Sustain Issue 2 (2023). <https://algorithmwatch.org/en/wp-content/uploads/2023/03/SustAIIn-magazine-issue-2.pdf>
- 150 “Artificial Intelligence Is Booming - So Is Its Carbon Footprint”, Josh Saul and Dina Bass, Bloomberg (2023). <https://www.bloomberg.com/news/articles/2023-03-09/how-much-energy-do-ai-and-chatgpt-use-no-one-knows-for-sure>
- 151 “Green AI”, Roy Schwartz, Jesse Dodge, Noah A. Smith and Oren Etzioni (2019). <https://arxiv.org/abs/1907.10597> (p.2)
- 152 “Tech companies underreport CO2 emissions”, Technical University of Munich, (2021). <https://www.sciencedaily.com/releases/2021/11/211118203514.htm>
- 153 “Green AI”, Roy Schwartz, Jesse Dodge, Noah A. Smith and Oren Etzioni (2019). <https://arxiv.org/pdf/1907.10597.pdf> (p.2)
- 154 “Generating Harms: Generative AI’s Impact & Paths Forward”, EPIC (2023). [EPIC-Generative-AI-White-Paper-May2023.pdf](https://www.epic-generative-ai-white-paper-may2023.pdf) (p.40)
- 155 “Green AI”, Roy Schwartz, Jesse Dodge, Noah A. Smith and Oren Etzioni (2019). <https://arxiv.org/pdf/1907.10597.pdf>
- 156 “Green AI”, Roy Schwartz, Jesse Dodge, Noah A. Smith and Oren Etzioni (2019). <https://arxiv.org/pdf/1907.10597.pdf>
- 157 “The mounting human and environmental costs of generative AI”, Sasha Luccion, ArsTechnica (2023). <https://arstechnica.com/gadgets/2023/04/generative-ai-is-cool-but-lets-not-forget-its-human-and-environmental-costs/>
- 158 “We’re getting a better idea of AI’s true carbon footprint”, Melissa Heikkilä, MIT Technology Review (2022) <https://www.technologyreview.com/2022/11/14/1063192/were-getting-a-better-idea-of-ais-true-carbon-footprint/>
- 159 “AI and the Challenge of Sustainability”, Sustain Issue 2 (2023). (p.16) <https://algorithmwatch.org/en/wp-content/uploads/2023/03/SustAIIn-magazine-issue-2.pdf>
- 160 “Green Intelligence: Why Data And AI Must Become More Sustainable”, Bernard Marr, Forbes (2023). <https://www.forbes.com/sites/bernardmarr/2023/03/22/green-intelligence-why-data-and-ai-must-become-more-sustainable/>
- 161 “AI machines aren’t ‘hallucinating’. But their makers are”, Naomi Klein, The Guardian (2023). <https://www.theguardian.com/mentisfree/2023/may/08/ai-machines-hallucinating-naomi-klein>
- 162 “Climate Change Impacts and Risks”, IPCC (2022). [https://www.ipcc.ch/report/ar6/wg2/downloads/outreach/IPCC\\_AR6\\_WGII\\_FactSheet\\_FoodAndWater.pdf](https://www.ipcc.ch/report/ar6/wg2/downloads/outreach/IPCC_AR6_WGII_FactSheet_FoodAndWater.pdf) and “Water – at the center of the climate crisis”, United Nations. <https://www.un.org/en/climate-change/science/climate-issues/water>
- 163 “Protect our people and future generations: Water and Climate Leaders call for urgent action”, World Meteorological Organization (2023). <https://public.wmo.int/en/media/press-release/protect-our-people-and-future-generations-water-and-climate-leaders-call-urgent>
- 164 OECD Environmental Outlook to 2050 (2012). [https://www.oecd-ilibrary.org/environment/oecd-environmental-outlook-to-2050\\_9789264122246-en](https://www.oecd-ilibrary.org/environment/oecd-environmental-outlook-to-2050_9789264122246-en)
- 165 “2022 Environmental Sustainability Report” Microsoft (p.28) <https://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RW14sJN>
- 166 “Making AI Less ‘Thirsty’: Uncovering and Addressing the Secret Water Footprint of AI Models”, Pengfei Li, Jianyi Yang, Mohammad A. Islam and Shaolei Ren (2023). <https://arxiv.org/pdf/2304.03271.pdf>
- 167 “‘Thirsty’ AI: Training ChatGPT Required Enough Water to Fill a Nuclear Reactor’s Cooling Tower, Study Finds”, Mack DeGeurin, Gizmodo (2023). <https://gizmodo.com/chatgpt-ai-water-185000-gallons-training-nuclear-1850324249>
- 168 “‘Thirsty’ AI: Training ChatGPT Required Enough Water to Fill a Nuclear Reactor’s Cooling Tower, Study Finds”, Mack DeGeurin, Gizmodo (2023). <https://gizmodo.com/chatgpt-ai-water-185000-gallons-training-nuclear-1850324249>
- 169 “Data centre water consumption”, David Mytton (2021). <https://www.nature.com/articles/s41545-021-00101-w>
- 170 “Microsoft will be carbon negative by 2030”, Microsoft (2020). Microsoft will be carbon negative by 2030 - The Official Microsoft Blog
- 171 “Google’s Carbon Offsets: Collaboration and Due Diligence”, Google (2011). <https://static>

- googleusercontent.com/media/www.google.com/no//green/pdfs/google-carbon-offsets.pdf, "2021 Sustainability Report", Meta <https://sustainability.fb.com/2021-sustainability-report/> and "2022 Environmental Sustainability Report", Microsoft. <https://www.microsoft.com/en-us/corporate-responsibility/sustainability/report>
- 172 Kaupa, C. (2022). "Peddling false solutions to worried consumers: The promotion of greenhouse gas 'offsetting' as a misleading commercial practice.", C. Kaupa, Journal of European Consumer and Market Law, 2022(4), 139-146. [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=4157810](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4157810)
- 173 "COP27: UN slams use of 'greenwashing' offsets ahead of direct abatement actions", Henry Edwardes-Evans, S&P (2022). <https://www.spglobal.com/commodityinsights/en/market-insights/latest-news/energy-transition/110922-cop27-un-slams-use-of-greenwashing-offsets-ahead-of-direct-abatement-actions>
- 174 "EU Parliament votes to clamp down on carbon neutral claims, early obsolescence", Valentina Romano, Euractiv (2023). <https://www.euractiv.com/section/energy-environment/news/eu-parliament-votes-to-clamp-down-carbon-neutral-claims-early-obsolescence/>
- 175 "Green AI", Roy Schwartz, Jesse Dodge, Noah A. Smith and Oren Etzioni (2019) <https://arxiv.org/pdf/1907.10597.pdf> (p. 2)
- 176 "Green AI", Roy Schwartz, Jesse Dodge, Noah A. Smith and Oren Etzioni (2019) <https://arxiv.org/pdf/1907.10597.pdf> (p. 7)
- 177 "Generating Harms: Generative AI's Impact & Paths Forward", EPIC (2023). [EPIC-Generative-AI-White-Paper-May2023.pdf](https://epic.org/new-epic-report-sheds-light-on-generative-ai-harms/) (p.41).
- 178 "OpenAI founder Sam Altman sees a big AI revolution within this decade", Matthias Bastian, The Decoder (2022). <https://the-decoder.com/openai-founder-sees-a-big-ai-revolution-within-this-decade/>
- 179 Ghostwork. <https://www.ghostwork.org>
- 180 "Exclusive: OpenAI Used Kenyan Workers on Less Than \$2 Per Hour to Make ChatGPT Less Toxic", Billy Perrigo, Time (2023). <https://time.com/6247678/openai-chatgpt-kenya-workers/>
- 181 Sama. <https://www.sama.com/why-sama/>
- 182 "AI and the Challenge of Sustainability", Sustain Issue 2 (2023). <https://algorithmwatch.org/en/wp-content/uploads/2023/03/Sustain-magazine-issue-2.pdf> (p. 12)
- 183 "AI and the Challenge of Sustainability", Sustain Issue 2 (2023). <https://algorithmwatch.org/en/wp-content/uploads/2023/03/Sustain-magazine-issue-2.pdf> (p. 11)
- 184 "Automation Shouldn't Always be Automatic: Making Artificial Intelligence Work for Workers and the World", Daron Acemoglu, OECD Forum (2020). <https://www.oecd-forum.org/posts/automation-shouldn-t-always-be-automatic-making-artificial-intelligence-work-for-workers-and-the-world>
- 185 "The New Generation of A.I. Apps Could Make Writers and Artists Obsolete", Nick Bilton, Vanity Fair (2022). <https://www.vanityfair.com/news/2022/06/the-new-generation-of-ai-apps-could-make-writers-and-artists-obsolete>
- 186 "Helpline workers for the National Eating Disorder Association say they are being replaced by AI", Britney Nguyen, Insider (2023). <https://www.businessinsider.com/eating-disorders-nonprofit-reportedly-fired-humans-offer-ai-chatbot-2023-5>
- 187 "Artists stage mass protest against AI-generated artwork on ArtStation", Benj Edwards, ArsTechnica (2022). <https://arstechnica.com/information-technology/2022/12/artstation-artists-stage-mass-protest-against-ai-generated-artwork/>
- 188 "The artist is dominating AI-generated art, and he's not happy about it. Melissa Heikkilä, MIT Technology Review (2022). <https://www.technologyreview.com/2022/09/16/1059598/this-artist-is-dominating-ai-generated-art-and-hes-not-happy-about-it/>
- 189 "We've filed a lawsuit challenging Stable Diffusion, a 21st century collage tool that violates the rights of artists", Stable Diffusion litigation (2023). <https://stablediffusionlitigation.com>
- 190 "Artists can now opt out of generative AI. It's not enough", Sayash Kapoor and Arvind Narayanan, AI Snake Oil (2023). <https://aisnakeoil.substack.com/p/artists-can-now-opt-out-of-generative>
- 191 "The Red Hand Files 218", Nick Cave (2023). <https://www.theredhandfiles.com/chat-gpt-what-do-you-think/>
- 192 "Artificially Yours: Who Owns Rights in AI-Generated Art?", Seyfarth Shaw (2023). <https://www.lexology.com/library/detail.aspx?q=640df36a-a63c-4936-82bd-579e4b-54ca00>
- 193 "Generating Harms: Generative AI's Impact & Paths Forward", EPIC (2023). <https://epic.org/new-epic-report-sheds-light-on-generative-ai-harms/>
- 194 Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation).
- 195 Art. 4(2) GDPR.
- 196 Art. 4(1) GDPR.
- 197 Art. 3 GDPR.
- 198 Art. 4(7) GDPR.
- 199 Art. 4(8) GDPR.
- 200 The GDPR does not apply to "natural persons' processing of data in the course of a purely personal or household activity", cf. art. 3. Generation of output relating to a natural person as part of purely personal activity, where the output is not shared online, is therefore not necessarily regulated by the GDPR. The GDPR can still apply to the system owner.
- 201 "Extracting Training Data from Large Language Models", Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Dawn Song, Katherine Lee, Adam Roberts, Tom Brown, Úlfar Erlingsson, Alina Oprea and Colin Raffel (2021). <https://arxiv.org/pdf/2012.07805.pdf>
- 202 "Algorithms that remember: Model inversion attacks and data protection law", Michael Veale, Reuben Binns and Lilian Edwards (2018). <https://royalsocietypublishing.org/doi/full/10.1098/rsta.2018.0083>
- 203 Art. 6 GDPR.
- 204 Art. 9(2) GDPR.
- 205 "ChatGPT: OpenAI reinstates service in Italy with enhanced transparency and rights

- for european users and non-users”, Italian DPA (2023). <https://www.garanteprivacy.it/home/docweb/-/docweb-display/docweb/9881490>
- 206 “OpenAI - Privacy policy”, (April 27, 2023). <https://openai.com/policies/privacy-policy>
- 207 “Google Bard hits over 180 countries and territories—none are in the EU”, Sharon Harding, ArsTechnica (2023). <https://arstechnica.com/gadgets/2023/05/google-bard-hits-over-180-countries-and-territories-none-are-in-the-eu/>
- 208 “More Penguins Than Europeans Can Use Google Bard”, Morgan Meaker and Matt Burgess, Wired (2023). <https://www.wired.co.uk/article/google-bard-european-union>
- 209 Bard FAQ (Accessed on June 1st, 2023). <https://bard.google.com/faq?hl=en>. Google Privacy Policy (December 15, 2022). <https://policies.google.com/privacy>
- 210 Art. 25 GDPR.
- 211 Art. 5(1)(c) GDPR.
- 212 Art. 5(1)(b) GDPR.
- 213 See for instance “How should we assess security and data minimisation in AI?”, ICO, <https://ico.org.uk/for-organisations/guide-to-data-protection/key-dp-themes/guidance-on-ai-and-data-protection/how-should-we-assess-security-and-data-minimisation-in-ai/>, laying out some guidance on this subject.
- 214 Art. 17 GDPR.
- 215 Art. 16 GDPR.
- 216 Art. 21 GDPR.
- 217 OpenAI Privacy policy 2023 (last update April 27, 2023). <https://openai.com/policies/privacy-policy>
- 218 See however requirements of information for data subjects in art. 12–14 GDPR.
- 219 “OpenAI’s hunger for data is coming back to bite it”. Melissa Heikkilä, MIT Technology Review (2023)(note that Open AI has not shared the size of the training data set for GPT-4.) <https://www.technologyreview.com/2023/04/19/1071789/openais-hunger-for-data-is-coming-back-to-bite-it/>
- 220 “Everyone wants to do the model work, not the data work”: Data Cascades in High-Stakes”, AI Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh and Lora Aroyo (2021). <https://storage.googleapis.com/pub-tools-public-publication-data/pdf/0d556e45afc54afeb2eb6b51a9b-c1827b9961ff4.pdf>
- 221 “Artificial intelligence: Stop to ChatGPT by the Italian SA” (2023). <https://www.gpdp.it/web/guest/home/docweb/-/docweb-display/docweb/9870847>
- 222 OpenAI Privacy policy 2023 (last update April 27, 2023). <https://openai.com/policies/privacy-policy>.
- 223 Umbrella term for AI systems designed to perform a wide range of tasks across different domains.
- 224 “ChatGPT: OpenAI reinstates service in Italy with enhanced transparency and rights for European users and non-users”, Italian DPA (2023). <https://www.garanteprivacy.it/home/docweb/-/docweb-display/docweb/9881490>
- 225 Art. 6(1)(f) GDPR.
- 226 Art. 6 and 5(2) GDPR.
- 227 “The Long and Winding Road: Two years of the GDPR: A cross-border data protection enforcement case from a consumer perspective”, BEUC (2020). [https://www.beuc.eu/sites/default/files/publications/beuc-x-2020-074\\_two\\_years\\_of\\_the\\_gdpr\\_a\\_cross-border\\_data\\_protection\\_enforcement\\_case\\_from\\_a\\_consumer\\_perspective.pdf](https://www.beuc.eu/sites/default/files/publications/beuc-x-2020-074_two_years_of_the_gdpr_a_cross-border_data_protection_enforcement_case_from_a_consumer_perspective.pdf)
- 228 “EDPB resolves dispute on transfers by Meta and creates task force on Chat GPT”, European Data Protection Board (2023). [https://edpb.europa.eu/news/news/2023/edpb-resolves-dispute-transfers-meta-and-creates-task-force-chat-gpt\\_en](https://edpb.europa.eu/news/news/2023/edpb-resolves-dispute-transfers-meta-and-creates-task-force-chat-gpt_en)
- 229 “Artificial intelligence: the action plan of the CNIL” (2023). <https://www.cnil.fr/en/artificial-intelligence-action-plan-cnil>
- 230 “European privacy watchdog creates ChatGPT task force”, Toby Sterling, Reuters, (2023). <https://www.reuters.com/technology/european-data-protection-board-discussing-ai-policy-thursday-meeting-2023-04-13/>
- 231 Directive 2005/29/EC of the European Parliament and of the Council of 11 May 2005 concerning unfair business-to-consumer commercial practices in the internal market and amending Council Directive 84/450/EEC, Directives 97/7/EC, 98/27/EC and 2002/65/EC of the European Parliament and of the Council and Regulation (EC) No 2006/2004 of the European Parliament and of the Council (Unfair Commercial Practices Directive).
- 232 Art. 2(b) UCPD.
- 233 Art. 4(d) UCPD.
- 234 Art. 6–7 UCPD.
- 235 Art. 8–9 UCPD.
- 236 “Guidance on the interpretation and application of Directive 2005/29/EC of the European Parliament and of the Council concerning unfair business-to-consumer commercial practices in the internal market” (2021). [https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52021XC1229\(05\)&from=EN](https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52021XC1229(05)&from=EN) (p. 12)
- 237 “Guidance on the interpretation and application of Directive 2005/29/EC of the European Parliament and of the Council concerning unfair business-to-consumer commercial practices in the internal market” (2021). [https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52021XC1229\(05\)&from=EN](https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52021XC1229(05)&from=EN) (p. 37)
- 238 “EU Consumer Protection 2.0 Structural asymmetries in digital consumer markets”, BEUC (2021). [https://www.beuc.eu/sites/default/files/publications/beuc-x-2021-018\\_eu\\_consumer\\_protection\\_2.0.pdf](https://www.beuc.eu/sites/default/files/publications/beuc-x-2021-018_eu_consumer_protection_2.0.pdf)
- 239 “That was fast! Microsoft slips ads into AI-powered Bing Chat”, Devin Coldewey, Techcrunch (2023). <https://techcrunch.com/2023/03/29/that-was-fast-microsoft-slips-ads-into-ai-powered-bing-chat/>
- 240 A European consumer umbrella organization of which the Norwegian Consumer Council is a member.
- 241 “Call for action to open an inquiry on generative AI systems to address risks and harms for consumers”, BEUC (2023). [https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023-045\\_Call\\_for\\_action\\_CPC\\_authorities\\_Generative\\_AI\\_systems.pdf](https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023-045_Call_for_action_CPC_authorities_Generative_AI_systems.pdf)
- 242 “Towards European Digital Fairness - BEUC framing response paper for the REFIT consultation”, BEUC (2023). <https://www.>



- [beuc.eu/sites/default/files/publications/BEUC-X-2023-020\\_Consultation\\_paper\\_RE-FIT\\_consumer\\_law\\_digital\\_fairness.pdf](https://beuc.eu/sites/default/files/publications/BEUC-X-2023-020_Consultation_paper_RE-FIT_consumer_law_digital_fairness.pdf)
- 243 “Digital fairness – fitness check on EU consumer law”, European Commission. [https://ec.europa.eu/info/law/better-regulation/have-your-say/initiatives/13413-Digital-fairness-fitness-check-on-EU-consumer-law\\_en](https://ec.europa.eu/info/law/better-regulation/have-your-say/initiatives/13413-Digital-fairness-fitness-check-on-EU-consumer-law_en)
- 244 Federal Trade Commission Act. <https://www.ftc.gov/legal-library/browse/statutes/federal-trade-commission-act>
- 245 “A Brief Overview of the Federal Trade Commission’s Investigative, Law Enforcement, and Rulemaking Authority”, FTC (2021). <https://www.ftc.gov/about-ftc/mission/enforcement-authority>
- 246 “Aiming for truth, fairness, and equity in your company’s use of AI”, FTC (2021). <https://www.ftc.gov/business-guidance/blog/2021/04/aiming-truth-fairness-equity-your-companys-use-ai>
- 247 “Keep your AI claims in check”, FTC (2023). <https://www.ftc.gov/business-guidance/blog/2023/02/keep-your-ai-claims-check>
- 248 “Model Destruction – The FTC’s Powerful New AI and Privacy Enforcement Tool”, Avi Gesser, Paul D. Rubin and Aanna Gresse, Debevoise Data Blog (2022). <https://www.debevoisedatablog.com/2022/03/22/model-destruction-the-ftcs-powerful-new-ai-enforcement-tool/>
- 249 Center for AI and Digital Policy. <https://www.caidp.org/cases/openai/>
- 250 “Joint statement on enforcement efforts against discrimination and bias in automated systems”, EEOC, CRT, FTC, CFPB-AI. [https://www.ftc.gov/system/files/ftc\\_gov/pdf/EEOC-CRT-FTC-CFPB-AI-Joint-Statement%28final%29.pdf](https://www.ftc.gov/system/files/ftc_gov/pdf/EEOC-CRT-FTC-CFPB-AI-Joint-Statement%28final%29.pdf)
- 251 “Biden–Harris Administration Announces New Actions to Promote Responsible AI Innovation that Protects Americans’ Rights and Safety”, The White House (2023). <https://www.whitehouse.gov/briefing-room/statements-releases/2023/05/04/fact-sheet-biden-harris-administration-announces-new-actions-to-promote-responsible-ai-innovation-that-protects-americans-rights-and-safety/>
- [cans-rights-and-safety/](https://cans-rights-and-safety/)
- 252 Directive 2001/95/EC of the European Parliament and of the Council of 3 December 2001 on general product safety.
- 253 Regulation (EU) 2023/988 of the European Parliament and of the Council of 10 May 2023 on general product safety, amending Regulation (EU) No 1025/2012 of the European Parliament and of the Council and Directive (EU) 2020/1828 of the European Parliament and the Council, and repealing Directive 2001/95/EC of the European Parliament and of the Council and Council Directive 87/357/EEC.
- 254 Art. 1(2) GPSD.
- 255 Art. 3(1) GPSD.
- 256 “Opinion of the sub-group on artificial Intelligence (ai), connected products and other new challenges in product safety to the consumer safety network”, European Commission (2021). [https://ec.europa.eu/safety/consumers/consumers\\_safety\\_gate/home/documents/Subgroup\\_opinion\\_final\\_format.pdf](https://ec.europa.eu/safety/consumers/consumers_safety_gate/home/documents/Subgroup_opinion_final_format.pdf)
- 257 Art. 2(1)(b) GPSD.
- 258 “Opinion of the sub-group on artificial Intelligence (ai), connected products and other new challenges in product safety to the consumer safety network”, European Commission (2021). [https://ec.europa.eu/safety/consumers/consumers\\_safety\\_gate/home/documents/Subgroup\\_opinion\\_final\\_format.pdf](https://ec.europa.eu/safety/consumers/consumers_safety_gate/home/documents/Subgroup_opinion_final_format.pdf)
- 259 Art. 8(2) GPSD.
- 260 “Urgent all for action regarding generative AI systems and concerns related to their safety”, BEUC (2023). [https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023046\\_BEUC\\_concerns\\_over\\_AI\\_and\\_mental\\_health\\_%20Ms\\_Pinuccia\\_Contino.pdf](https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023046_BEUC_concerns_over_AI_and_mental_health_%20Ms_Pinuccia_Contino.pdf)
- 261 Cf. recital 19 GPSR.
- 262 Art. 6(1)(h) GPSR.
- 263 Art. 101 TFEU.
- 264 Art. 102 TFEU.
- 265 “Competition policy”. European Parliament (2023). [https://www.europarl.europa.eu/erpl-app-public/factsheets/pdf/en/FTU\\_2.6.12.pdf](https://www.europarl.europa.eu/erpl-app-public/factsheets/pdf/en/FTU_2.6.12.pdf)
- 266 Regulation (EU) 2022/2065 of the European Parliament and of the Council of 19 October 2022 on a Single Market For Digital Services and amending Directive 2000/31/EC (Digital Services Act).
- 267 Art. 3(1)(g) DSA.
- 268 “DSA: Very large online platforms and search engines”, European Commission. <https://digital-strategy.ec.europa.eu/en/policies/dsa-vlops>
- 269 “Understanding and Regulating ChatGPT, and Other Large Generative AI Models”, Philipp Hacker, Andreas Engel, Theresa List, Verfassungsblog (2023). <https://verfassungsblog.de/chatgpt/>
- 270 Proposal for a regulation of the European Parliament and of the Council laying down harmonised rules on artificial intelligence and amending certain union legislative acts, para. 5 (2021). [https://eurlex.europa.eu/resource.html?uri=cellar:e0649735-a372-11eb-9585-01aa75e-d71a1.0001.02/DOC\\_1&format=PDF](https://eurlex.europa.eu/resource.html?uri=cellar:e0649735-a372-11eb-9585-01aa75e-d71a1.0001.02/DOC_1&format=PDF)
- 271 Art. 2(1)(a) Draft AIA.
- 272 Art. 3(1) Draft AIA, cf. Annex 1.
- 273 Art. 5 Draft AIA.
- 274 Art. 6 Draft AIA.
- 275 See Draft AIA Title III.
- 276 Art. 17 Draft AIA.
- 277 Art. 9 Draft AIA.
- 278 Art. 10 Draft AIA.
- 279 Art. 15 Draft AIA.
- 280 Art. 11 Draft AIA.
- 281 Art. 52 Draft AIA.
- 282 Art. 69 Draft AIA.
- 283 “Regulation AI to protect the consumer – Position Paper on the AI Act”, BEUC (2021). [https://www.beuc.eu/sites/default/files/publications/beuc-x-2021-088\\_regulating-ai-to-protect-the-consumer.pdf](https://www.beuc.eu/sites/default/files/publications/beuc-x-2021-088_regulating-ai-to-protect-the-consumer.pdf)
- 284 Art. 5 Draft AIA.
- 285 “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 – General approach”, Council of the European Union (2022). <https://data.consilium.europa.eu/doc/document/ST-14954-2022-INIT/en/pdf>
- 286 Art. 3(1)(b) Council Position.
- 287 Art. 4b(1) Council Position.

- 288 Art. 4c(1) Council Position.
- 289 Art. 4c(2) Council Position.
- 290 “DRAFT Compromise Amendments on the Draft Report – Proposal for a regulation of the European Parliament and of the Council on harmonised rules on Artificial Intelligence (Artificial Intelligence Act) and amending certain Union Legislative Acts”, LIBE and IMCO committees of the European Parliament. (2023). [https://www.europarl.europa.eu/meetdocs/2014\\_2019/plmrep/COMMITTEES/CJ40/DV/2023/05-11/ConsolidatedCA\\_IMCOLIBE\\_AI\\_ACT\\_EN.pdf](https://www.europarl.europa.eu/meetdocs/2014_2019/plmrep/COMMITTEES/CJ40/DV/2023/05-11/ConsolidatedCA_IMCOLIBE_AI_ACT_EN.pdf)
- 291 Art. 68c Parliament Position.
- 292 Art. 68a Parliament Position.
- 293 Art. 68b Parliament Position.
- 294 Art. 68d Parliament Position.
- 295 Art. 3(1c) Parliament Position.
- 296 Art. 28b Parliament Position.
- 297 For the complete list of requirements, see art. 28b(2) Parliament Position.
- 298 Art. 28b(4).
- 299 “ChatGPT and the AI Act”, Natali Helberger and Nicholas Diakopoulos (2023). <https://policyreview.info/essay/chatgpt-and-ai-act>
- 300 “General Purpose AI Poses Serious Risks, Should Not Be Excluded From the EU’s AI Act”, Amba Kak and Sarah Myers West, AI Now Institute (2023). <https://ainowinstitute.org/publication/gpai-is-high-risk-should-not-be-excluded-from-eu-ai-act>
- 301 “The lobbying ghost in the machine”, Corporate Europe Observatory (2023). <https://www.corporateeurope.org/en/2023/02/lobbying-ghost-machine>
- 302 The exact timing differs between the different positions of the EU institutions. At the earliest the AIA will be fully applicable after 24 months after entering into force. This effectively means it may not be fully applicable until at the earliest April/May 2026, if there is a trilogue agreement by January 2024.
- 303 Case C-65/20, VI v Krone (2021), <https://curia.europa.eu/juris/document/document.jsf?jsessionid=7A0662FAD-49ED462BE89A81594FAF809?text=&docid=242561&pageIndex=0&doclang=EN>.
- 304 “Revision of the Product Liability Directive”, BEUC (2023). [https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023-023\\_Revision\\_of\\_the\\_product\\_liability\\_directive.pdf](https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023-023_Revision_of_the_product_liability_directive.pdf)
- 305 “Proposal for an AI Liability Directive”, BEUC (2023). [https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023-050\\_Proposal\\_for\\_an\\_AI\\_Liability\\_Directive.pdf](https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023-050_Proposal_for_an_AI_Liability_Directive.pdf)
- 306 “Proposal for an AI Liability Directive”, BEUC (2023). [https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023-050\\_Proposal\\_for\\_an\\_AI\\_Liability\\_Directive.pdf](https://www.beuc.eu/sites/default/files/publications/BEUC-X-2023-050_Proposal_for_an_AI_Liability_Directive.pdf)
- 307 “How to create, release, and share generative AI responsibly”, Melissa Heikkilä, MIT Technology Review (2023). <https://www.technologyreview.com/2023/02/27/1069166/how-to-create-release-and-share-generative-ai-responsibly/>
- 308 “Pause Giant AI Experiments: An Open Letter”, Future of Life (2023). <https://futureoflife.org/open-letter/pause-giant-ai-experiments/>
- 309 “OpenAI’s CEO confirms the company isn’t training GPT-5 and ‘won’t for some time’”, James Vincent, The Verge (2023). <https://www.theverge.com/2023/4/14/23683084/openai-gpt-5-rumors-training-sam-altman>
- 310 “EU, Google to develop voluntary AI pact ahead of new AI rules, EU’s Breton says”, Foo Yun Chee, Reuters (2023). <https://www.reuters.com/technology/eu-google-develop-voluntary-ai-pact-ahead-new-ai-rules-eu-breton-says-2023-05-24/>



**FOR MORE INFORMATION:**

Finn Lützow-Holm Myrstad, director of Digital Policy  
Norwegian Consumer Council  
E-mail: [finn.myrstad@forbrukerradet.no](mailto:finn.myrstad@forbrukerradet.no)  
[www.forbrukerradet.no/ai](http://www.forbrukerradet.no/ai)

