



2 AI content detection in the emerging information ecosystem: new 3 obligations for media and tech companies

4 Alistair Knott^{1,2} · Dino Pedreschi^{1,3} · Toshiya Jitsuzumi^{1,4} · Susan Leavy^{1,5} · David Eyers^{1,6} ·
5 Tapabrata Chakraborti^{1,7,8} · Andrew Trotman⁶ · Sundar Sundareswaran¹ · Ricardo Baeza-Yates^{1,9} ·
6 Przemyslaw Biecek^{1,10} · Adrian Weller^{1,11} · Paul D. Teal^{1,12} · Subhadip Basu^{1,13} · Mehmet Haklidir^{1,14} ·
7 Virginia Morini^{1,3} · Stuart Russell^{1,15} · Yoshua Bengio^{1,16,17}

8 Accepted: 5 August 2024
9 © The Author(s), under exclusive licence to Springer Nature B.V. 2024

10 Abstract

11 The world is about to be swamped by an unprecedented wave of AI-generated content. We need reliable ways of identifying
12 such content, to supplement the many existing social institutions that enable trust between people and organisations and
13 ensure social resilience. In this paper, we begin by highlighting an important new development: providers of AI content
14 generators have new obligations to *support* the creation of reliable detectors for the content they generate. These new obli-
15 gations arise mainly from the EU's newly-finalised AI Act, but they are enhanced by the US President's recent Executive
16 Order on AI, and by several considerations of self-interest. These new steps towards reliable detection mechanisms are by
17 no means a panacea—but we argue they will usher in a new adversarial landscape, in which reliable methods for identify-
18 ing AI-generated content are commonly available. In this landscape, many new questions arise for policymakers. Firstly, if
19 reliable AI-content detection mechanisms are available, *who should be required to use them?* And how should they be used?
20 We argue new duties arise for media companies, and for Web search companies, in the deployment of AI-content detectors.
21 Secondly, what broader regulation of the tech ecosystem will maximise the likelihood of reliable AI-content detectors? We
22 argue for a range of new duties, relating to provenance-authentication protocols, open-source AI generators, and support for
23 research and enforcement. Along the way, we consider how the production of AI-generated content relates to 'free expres-
24 sion', and discuss the important case of content that is generated jointly by humans and AIs.

25 **Keywords** Generative AI · AI-generated content · AI regulation

A1 Alistair Knott
A2 ali.knott@vuw.ac.nz

A3 ¹ Social Media Governance Project, Global Partnership on AI,
A4 Montreal, Canada

A5 ² School of Engineering and Computer Science, Victoria
A6 University of Wellington, Wellington, New Zealand

A7 ³ University of Pisa, Pisa, Italy

A8 ⁴ Chuo University, Tokyo, Japan

A9 ⁵ Insight SFI Research Centre for Data Analytics, School
A10 of Information and Communication, University College
A11 Dublin, Dublin, Ireland

A12 ⁶ School of Computing, University of Otago, Dunedin,
A13 New Zealand

A14 ⁷ Alan Turing Institute, London, United Kingdom

A15 ⁸ University College London, London, United Kingdom

⁹ Institute for Experiential AI, Northeastern University,
Silicon Valley, USA

¹⁰ Warsaw University of Technology, Warsaw, Poland

¹¹ Director of Research, Machine Learning, University
of Cambridge, Cambridge, United Kingdom

¹² Victoria University of Wellington, Wellington, New Zealand

¹³ Computer Science and Engineering Department, Jadavpur
University, Kolkata, India

¹⁴ Artificial Intelligence Institute, Tubitak Bilgem, Gebze,
Türkiye

¹⁵ Center for Human-Compatible AI, UC Berkeley, Berkeley,
USA

¹⁶ Mila - Quebec AI Institute, Montreal, Canada

¹⁷ University of Montreal, Montreal, Canada

A16
A17

A18

A19

A20

A21

A22

A23

A24

A25

A26

A27

A28

A29

26

Introduction

27 The Web, and the world beyond it, are about to be
 28 swamped by a wave of AI-generated content. AI text gen-
 29 eration systems, such as GPT-4 (OpenAI, 2023) Gemini
 30 (Google, 2024), Llama (Touvron et al., 2023), Falcon
 31 (UAE TII, 2023) and Mixtral (Jiang et al., 2024), are
 32 becoming widely used to produce textual content in a vari-
 33 ety of domains, such as news (Newsguard, 2024), business
 34 reviews (Berry, 2024), academia (Originality, 2024) and
 35 culture (Notopoulos, 2024), in an extensive range of lan-
 36 guages (see e.g. Fernandes, 2023). AI image generation
 37 systems, such as Dall-E (OpenAI, 2021) and MidJourney
 38 (Midjourney, Inc., 2022) are producing huge volumes of
 39 AI-generated content online (see e.g. Valyaeva, 2023),
 40 and are radically changing workflows for human graphic
 41 designers (see e.g. HackerNoon, 2023). Images seem likely
 42 soon to be followed by AI video generation, such as Sora
 43 (OpenAI, 2024).

44 The widespread adoption of AI content-generation tech-
 45 nologies brings many benefits (see Dell'Acqua et al., 2023;
 46 Candelon et al., 2023 for balanced reviews). However, this
 47 proliferation of AI-generated content also presents signifi-
 48 cant challenges. As AI generation systems improve, it will
 49 become increasingly difficult for human consumers of con-
 50 tent to accurately tell whether an item of content was pro-
 51 duced by a person or an AI system, or some combination
 52 of the two. This poses a brand new authentication problem:
 53 as the differences between AI-generated and human-gen-
 54 erated content decrease, it becomes intrinsically harder to
 55 adjudicate individual cases.

56 Why do we need to know whether an item was gener-
 57 ated by a person or an AI? Importantly, the reasons don't
 58 hinge on the *quality* of the content. Human-generated
 59 content and AI-generated content can both vary enor-
 60 mously in quality. In the right contexts, both humans and
 61 AIs can produce useful, truthful, informative content; in
 62 other contexts, both humans and AIs are capable of pro-
 63 ducing harmful, misleading, inaccurate content. The rea-
 64 sons rather hinge on the role of AI content generation as a
 65 *social practice*. Communication between humans through
 66 the creation of enduring content (text, images and other
 67 media) is fundamental to the ordering of our societies:
 68 human-generated content plays a central role in the crea-
 69 tion and enforcement of laws, in education and training, in
 70 the dissemination of news and opinion, in the organisa-
 71 tion of political debates and democratic processes, in the for-
 72 mation and transmission of culture. In all these contexts,
 73 societies have developed resilient institutions that allow
 74 citizens to have confidence in human-generated content:
 75 from educational providers that certify individuals as
 76 reputable content providers in specific domains, to laws

77 governing the broadcasting of content and the function-
 78 ing of political debates, to conventions about the rule of
 79 law. AI-generated content *escapes* many of our existing
 80 institutions.

81 AI content generation escapes existing institutions in
 82 two main ways. Firstly, it lets people *deliver content they*
 83 *didn't produce*, and maybe don't even understand. In many
 84 cases they may not even have seen or read it. In educational
 85 settings, this undermines traditional assessment practices,
 86 and disrupts current accreditation systems. It also appears
 87 to be impacting academic review processes (see Liang
 88 et al., 2024). In the professional world, AI content genera-
 89 tion undermines the processes through which people and
 90 organisations acquire reputations for reliable work. In all
 91 these cases, AI threatens breakdowns of social trust. Sec-
 92 ondly, AI lets people *proliferate content*. A single person can
 93 produce vastly more content than before, including content
 94 carefully tailored to specific audiences. This allows indi-
 95 viduals to exert new and unprecedented influences on public
 96 discussions. The new influences in political discussions are
 97 particularly concerning: the recent deepfake of Joe Biden's
 98 voice (NBC, 2024) provides a taste of what is now possible.
 99 Organisations can similarly increase their capacity to pro-
 100 duce content with generative AI, so organisations also have
 101 new powers of influence on public discussions. The fact that
 102 public discussions increasingly happen online amplifies the
 103 effects of these new abilities to proliferate content, and to
 104 add coherently to existing content. And AI-generated con-
 105 tent is known to have effects in changing consumers' senti-
 106 ment; see for instance Jakesch et al. (2023).

107 In short, AI content generation systems can pose serious
 108 threats to social stability, and especially to political stabil-
 109 ity. 2024 will see democratic elections taking place across
 110 the globe, so these threats are immediate. To counter these
 111 threats, we need to *extend* the institutions that currently
 112 govern content creation, to make provisions for generative
 113 AI. The crucial extension is to provide methods of *reliably*
 114 *identifying AI-generated content*, and reliably distin-
 115 guishing it from human-generated content. Finding such meth-
 116 ods involves tackling several related questions, which bear
 117 on technical and legal mechanisms, but also on economics
 118 and company incentives, and on the operation of the open-
 119 source ecosystem. In two recent papers (GPAI, 2023; Knott
 120 et al., 2023) we reviewed these questions, and argued that
 121 the best way to obtain reliable mechanisms for detecting
 122 AI-generated content is to place responsibility for the pro-
 123 vision of these mechanisms with the organisations (princi-
 124 pally companies) that build and deploy generative AI tools.
 125 Specifically, we proposed that any agency that creates an AI
 126 content generator must be required to *demonstrate a reliable*
 127 *detection mechanism* for the content that generator produces,
 128 as a *condition of its use by the public*—and to make the
 129 detection mechanism publicly available (as a closed-source

130 tool) on its release. See GPAI (2023); Knott et al. (2023)
 131 for details of this proposal. (We will discuss what counts as
 132 'reliable' later in the paper.)

133 Our proposal, along with some allied efforts we will dis-
 134 cuss, had good traction with policymakers in the EU and
 135 the US: it was influential in shaping some new legal and
 136 organisational directives for generative AI providers. In the
 137 second section of this paper, we will review these new direc-
 138 tives. In the third section, we take stock of the new land-
 139 scape for AI-generated content detection which these new
 140 directives set up. The directives are certainly not a panacea.
 141 Instead, we argue they set the stage for an ongoing 'arms
 142 race', between *providers* of AI content detectors (both inside
 143 and outside generator companies) and actors who seek to
 144 *evade* detection. In this new landscape, we expect that reli-
 145 able methods for discriminating between AI-generated and
 146 natural or human-generated content will sometimes—per-
 147 haps often—be available.

148 This analysis prompts two new sets of questions for poli-
 149 cymakers. Firstly, if reliable methods exist for identifying
 150 AI-generated content, *who should use these methods?* And
 151 *how should they be used?* We consider these questions in the
 152 fourth section of the paper, and conclude with some recom-
 153 mendations about new rules for media companies, and per-
 154 haps for Web search companies. Secondly, what policy steps
 155 can be taken to *intervene* in the arms race between providers
 156 and evaders of AI-content identification systems, to ensure
 157 that reliable identification methods are widely and frequently
 158 available? We consider this question in the fifth section of
 159 the paper, and conclude with recommendations about several
 160 aspects of the broader information ecosystem.

161 **New imperatives on AI providers 162 regarding AI-generated content 163 identification**

164 **Obligations imposed by the EU's AI Act**

165 The EU's AI Act, whose final text has recently been agreed
 166 (see e.g. EU/FLI, 2024), explicitly recognises the potential
 167 of AI-generated content to destabilise society, and the role
 168 AI providers should play to prevent this. As stated in Recital
 169 70a:

170 A variety of AI systems can generate large quantities
 171 of synthetic content that becomes increasingly hard
 172 for humans to distinguish from human-generated and
 173 authentic content. The wide availability and increasing
 174 capabilities of those systems have a significant impact
 175 on the integrity and trust in the information ecosystem
 176 (...) In the light of those impacts, (...) it is appropri-
 177 ate to require providers of those systems to embed

178 technical solutions that enable marking in a machine
 179 readable format and detection that the output has been
 180 generated or manipulated by an AI system and not a
 181 human. Such techniques and methods should be suf-
 182 ficiently reliable, interoperable, effective and robust as
 183 far as this is technically feasible, taking into account
 184 available techniques or a combination of such tech-
 185 niques, such as watermarks, metadata identifications,
 186 cryptographic methods for proving provenance and
 187 authenticity of content, logging methods (...)

188 The Act imposes some clear obligations on providers,
 189 which are stated in Article 52.1(a):

190 Providers of AI systems, including [General-Purpose
 191 AI] systems, generating synthetic audio, image, video
 192 or text content, shall ensure the outputs of the AI sys-
 193 tem are marked in a machine-readable format and
 194 detectable as artificially generated or manipulated.
 195 Providers shall ensure their technical solutions are
 196 effective, interoperable, robust and reliable as far as
 197 this is technically feasible, taking into account spe-
 198 cificities and limitations of different types of content,
 199 costs of implementation and the generally acknowl-
 200 edged state-of-the-art, as may be reflected in relevant
 201 technical standards. This obligation shall not apply to
 202 the extent the AI systems perform an assistive function
 203 for standard editing or do not substantially alter the
 204 input data provided by the deployer or the semantics
 205 thereof, or where authorised by law to detect, prevent,
 206 investigate and prosecute criminal offences.

207 Four comments are useful here. Firstly, obligations about
 208 content detection are only imposed for AI systems that gen-
 209 erate substantially new content; systems that make minor
 210 changes to existing content are sensibly exempted.

211 Secondly, obligations are subject to considerations of
 212 cost and technical feasibility, and reference is made to cer-
 213 tain types of content where technical challenges are higher.
 214 (Watermarking is more challenging for textual content than
 215 for images, for instance, as discussed by Srinivasan, 2024.)

216 Thirdly, note that the EU directive only refers to specific
 217 detection mechanisms (like watermarking) as *examples* of
 218 mechanisms that could function to support detection. The
 219 directive itself is rightly more general, accommodating the
 220 possibility that detection mechanisms may need to change
 221 as technology advances. Note that Recital 70a usefully refers
 222 to 'logging methods', which are a promising alternative to
 223 watermarking, but have received less attention. In these
 224 methods, the provider of the AI generator keeps a private log
 225 of content it generates (see Krishna et al., 2023 for the origi-
 226 nal proposal). A detector for the AI-generated content can
 227 then be implemented very simply as a plagiarism detector
 228 for content in this log, using mature Information Retrieval

229 technology. Further discussion of possible detection mechanisms, along with their pros and cons, is provided in Knott
 230 et al. (2023).¹

232 Finally, the mechanisms foreseen for detection include
 233 mechanisms for *proving provenance* (at least in Recital
 234 70a). The issue of provenance detection is broader than
 235 the issue of AI-generated content detection: several groups
 236 have suggested that the problems of AI-generated content
 237 are best addressed by a broader protocol that allows
 238 *human-generated content* to be positively authenticated.
 239 That proposal is particularly associated with the Content
 240 Authenticity Initiative and Project Origin, whose efforts are
 241 unified in the C2PA standard. The aim is that this standard
 242 is adopted throughout the ecosystem for capturing or
 243 generating, transforming, transmitting and viewing content.
 244 The standard could be adopted by camera manufacturers,
 245 for instance, to embed information about when and where a
 246 photo or video was recorded, or by broadcasters and other
 247 media organisations, to retain this embedded information. Of
 248 course these wider obligations don't belong in a piece of legis-
 249 lation about AI—but it is useful that the AI Act mentions
 250 the provenance-authentication proposal in a recital accom-
 251 panying obligations on generative AI providers to support
 252 detection. We will consider broader legislation supporting
 253 provenance-authentication later in this paper. (For now,,
 254 we will use the term ‘content *identification*’ to encompass
 255 both focussed AI-content detection and broader provenance-
 256 tracking schemes.)

257 **Guidance from Biden's executive order on AI**

258 In the US, President Biden issued an Executive Order ‘on
 259 the Safe, Secure, and Trustworthy Development and Use
 260 of AI’ in October last year. This order followed a Senate
 261 Judiciary Enquiry on ‘Oversight of AI’, at which two of our
 262 co-authors (Yoshua Bengio and Stuart Russell) gave evi-
 263 dence (alongside Dario Amodei from Anthropic). Much of
 264 the conversation at this Enquiry was about AI-generated
 265 content identification—and again, the methods discussed
 266 included mechanisms focussed specifically on AI-generated
 267 content detection tools, and broader protocols for tracking
 268 the provenance of all content, whether human- or AI-generated.
 269 The Executive Order aims to strengthen public trust in
 270 the authenticity of government communications, and more
 271 generally, to tackle disinformation. To these ends, it asks for
 272 a review of work on AI content detection in Sect. 4.5.(a):

1 It is worth noting that *combinations of different* detection mechanisms are likely to be particularly effective in delivering reliable detectors. Ensemble techniques for classification are likely to be beneficial here, just as they are elsewhere in machine learning (Zhou et al., 2014). We feel such ensemble methods are not yet widely enough discussed in relation to AI-content detection.

273 the Secretary of Commerce (...) shall submit a report
 274 (...) identifying the existing standards, tools, methods,
 275 and practices, as well as the potential development
 276 of further science-backed standards and techniques,
 277 for (...) (ii) labeling synthetic content, such as using
 278 watermarking; (iii) detecting synthetic content (...)

279 and for guidance about both detection and provenance-
 280 authentication in Sect. 4.5.(b):

281 the Secretary of Commerce, in coordination with the
 282 Director of OMB [the Office of Management and
 283 Budget], shall develop guidance regarding the existing
 284 tools and practices for digital content authentication
 285 and synthetic content detection measures (...)

286 In Sect. 10.1.(b) (viii)(c), the Director of OMB is addi-
 287 tionally tasked with making.

288 recommendations to [executive departments and]
 289 agencies regarding (...) reasonable steps to watermark
 290 or otherwise label output from generative AI[.]

291 These actions don't impose legal obligations on com-
 292 panies, but they directly impact government procurement
 293 processes, and create expectations that may have impacts
 294 in civil lawsuits.

295 **Obligations arising from the self-interest of AI 296 providers**

297 Alongside external guidance from policymakers, some new
 298 research findings give generative AI providers strong incen-
 299 tives of their own to support the detection of AI-generated
 300 content. If an AI generator re-trains on the content it pro-
 301 duced itself, its quality deteriorates substantially: a phe-
 302 nomenon termed ‘model collapse’, first reported by Shumailov et al.
 303 (2023) and now receiving much attention (see e.g. Dohmatob et al., 2024a, 2024b). AI providers therefore have good
 304 reason to exclude AI-generated content from their training
 305 sets—and thus have good incentives to be able to identify
 306 such content reliably. Note that providers also have separate
 307 (positive) incentives to identify text from their own genera-
 308 tors, to gauge uptake of their systems, which is a com-
 309 mercially important measure of performance.

310 Of course, companies may not want to impose a blanket
 311 ban on AI-generated training items. There are several sit-
 312 uations where AI-generated training items can help address
 313 issues in the dataset, such as data scarcity and bias (see e.g.
 314 de Wilde et al., 2024), and to augment data quality (for
 315 instance by removing noise, normalising, or increasing
 316 resolution). These *directed* uses of AI-content can be very
 317 beneficial; model collapse arises when the model’s training
 318 set is *indiscriminately* extended with AI-content.

320 **Summary**

321 Taken together, the new legal requirements about to be
 322 imposed in the EU, the recent guidance from Biden's Executive
 323 Order, and recently-recognised considerations of corporate
 324 self-interest allow us to confidently anticipate new
 325 initiatives from companies in support of AI content detection.
 326 The very recent 'Munich accord' in which 20 of the
 327 leading tech companies pledge to 'work together to detect
 328 and counter harmful AI content' in this year's elections
 329 (Munich, 2024) is some testament to this. The implementation
 330 and enforcement of these new initiatives will of course
 331 be challenging: we will review the main challenges in the
 332 next section.

333 Of the obligations discussed in the current section, we
 334 should note that by far the most stringent are those imposed
 335 by the EU, which require providers operating in the EU
 336 market to support detection mechanisms. As an aside, the
 337 largest AI generator companies, which will be centre stage
 338 for EU regulators, may sometimes deploy the same generators
 339 beyond the EU as within it. For detection methods that
 340 are built into generators, this may mean that EU-mandated
 341 support for detection will naturally extend to jurisdictions
 342 outside the EU. We feel there are good prospects for a 'Brus-
 343 sels effect' in this area, as has been found in other areas of
 344 EU tech legislation (Bradford, 2020).

345 **The new adversarial landscape for AI 346 content identification**

347 In the previous section, we reviewed a range of new obligations
 348 on providers of AI generators, to support reliable methods
 349 for identifying the content their systems generate. These
 350 obligations should prompt great improvements in the quality
 351 of methods for identifying AI-generated content—especially
 352 given the 'Brussels effect' we anticipated above. If the big
 353 AI companies fully engage with the goal of creating reliable
 354 detectors, we can expect reliable detectors to emerge, which
 355 are serviceable in the EU and some way beyond. Note that
 356 reliable detectors can also be expected to emerge from time
 357 to time even without support from providers. For instance,
 358 the recent methods for detecting images generated by stable
 359 diffusion (see Wang et al., 2023; Zhang and Xu, 2023) are
 360 impressively reliable; recent zero-shot methods for detecting
 361 LLM-generated text (e.g. Hans et al., 2024; Su et al.,
 362 2023) also show some promise, as do models fine-tuned for
 363 specific domains (see e.g. Veselovsky et al., 2023).

364 Of course, these are just the opening moves in a new, and
 365 doubtless ongoing, adversarial process. Any reliable method
 366 for AI-content detection, whether supported by providers,
 367 or developed externally, will trigger responses from actors

368 who wish to *evade* detection. For detectors that rely on finding
 369 differences between AI-generated and 'natural' content,
 370 there is an obvious point of attack: as noted by Májovský
 371 et al. (2024), any identified difference can immediately serve
 372 as an error term to train a new generator that eliminates
 373 exactly that difference. Detectors can also be attacked by
 374 *manipulating AI-generated content*, so it evades detection.
 375 For instance, changing some of the words in a generated text
 376 can destroy watermarks added by a generator (see e.g. Sada-
 377 sivan et al., 2023). Automated tools for modifying images, or
 378 paraphrasing texts, can likewise defeat detectors.² An early
 379 summary of this adversarial landscape is given by Crothers
 380 et al., (2023); a more recent summary is provided in a recent
 381 report by the Forum for Information and Democracy (FID,
 382 2024 Ch1 Sect. 1.5).

383 Fortunately, the drafters of the AI Act have anticipated
 384 these adversarial responses. Article 52.1(a) requires that
 385 AI company support for detection mechanisms be adequate
 386 given '*the generally acknowledged state-of-the-art*', which
 387 should certainly be understood to include known adversarial
 388 techniques. The AI Act can therefore be seen as defining
 389 providers' obligations in the 'arms race' which is now get-
 390 ting underway between the creators of detector tools (both
 391 within generator companies and beyond) and those attempt-
 392 ing to evade detection. The picture is complicated by actors
 393 who are reluctant to comply with existing rules, or unaware
 394 of these rules. The open-source software ecosystem poses
 395 some special challenges, both for enforcement of rules and
 396 in providing platforms for exploring adversarial strategies
 397 (as we will discuss further below). Whenever current meth-
 398 ods for identifying AI content are defeated, this will prompt
 399 the development of improved methods. It may be at certain
 400 points that the evaders have the upper hand, and AI provid-
 401 ers must work to find new ways of meeting their obligations.
 402 (Again, the AI Act provides for this contingency, by making
 403 providers' obligations subject to 'technical feasibility'.)
 404 Of course, arms races are nothing new for tech companies:
 405 Google has an ongoing battle with search engine optimisers
 406 (see e.g. Davis, 2006); social media companies have similar
 407 battles with purveyors of harmful content (see e.g. Founta
 408 et al., 2019). But it is useful to clearly identify the battle that
 409 is newly emerging between providers of AI-content detectors
 410 and those aiming to evade detection.

411 In this new adversarial and dynamic context, we foresee
 412 several new questions for policymakers. Firstly, if reliable
 413 methods for identifying AI-generated content are available
 414 at a given moment, *who should make use of them?* And *how*
 415 *should they be properly used?* We will consider those ques-
 416 tions in the next section. Secondly, what can policymakers

² Logging methods appear more resilient to paraphrase attacks, how-
 ever, as reported by Krishna et al. (2023).

417 do to *stack* the arms race in favour of reliable detection
 418 mechanisms? We will consider that question in the section
 419 after that.

420 **When reliable AI-content identification
 421 methods become available, who should
 422 make use of them?**

423 In this section, we will consider a scenario where reliable
 424 methods for identifying AI-generated content are available.
 425 In this scenario, policymakers need to determine *who should*
 426 *make use of* these reliable methods, and what constitutes
 427 their proper use.

428 A key consideration for policymakers relates to the *incentives*
 429 that ensure the proper use of identification methods
 430 within the information ecosystem. We begin by arguing that
 431 many organisations in society will naturally adopt reliable
 432 methods as they become available, as an organic extension
 433 of their existing mechanisms for maintaining reputation
 434 and trustworthiness amongst those they interact with. We
 435 then consider the case of media organisations. We argue
 436 that some of these organisations aren't naturally motivated
 437 to adopt systematic AI-generated content identification poli-
 438 cies, and hence should be required to do so by law. We con-
 439 sider various ways media companies could moderate the
 440 AI-generated content they detect. We conclude by surveying
 441 the many risks that arise in the process of identifying and
 442 moderating AI-generated content, and consider how policies
 443 can balance these against the risks arising from proliferation
 444 of AI content.

445 **Free-market incentives to use reliable AI-content
 446 identification methods**

447 As we discussed in the first section, AI content generation
 448 lets people deliver work that is not their own, that they may
 449 have had minimal involvement in, and may not have thor-
 450oughly checked. (We are thinking particularly here of AI-
 451 generated *text*, where the process of checking or vetting can
 452 require a considerable amount of human work.) This creates
 453 potential accountability gaps in any organisation where con-
 454 tent is to be produced. For instance, in educational institu-
 455 tions, students can deliver work they didn't produce or don't
 456 fully understand, which threatens the accreditations these
 457 institutions provide. In the professional world, workers can
 458 likewise deliver content they didn't produce, and can't fully
 459 vouch for, which threatens to undermine the credibility of
 460 individuals, and more importantly of whole organisations.

461 These problems are exacerbated by the tendency of AI
 462 generators to 'hallucinate' (see e.g. Rawte et al., 2023). This
 463 tendency can be mitigated in various ways (see e.g. Tonmoy
 464 et al., 2024), but it is still an inherent feature in systems

465 that are optimised on the surface form of training items,
 466 rather than on more abstract measures of meaning. But even
 467 disregarding hallucinations, there is a deeper problem: AI
 468 content generation potentially lets human providers 'fall out
 469 of the loop' in a professional relationship (see e.g. Zerilli
 470 et al., 2019). There is no guarantee that services are being
 471 provided by the people or companies who are contracted to
 472 do the work. Again, this leads to a huge accountability gap.

473 If reliable ways of identifying AI-generated content
 474 become available, we believe the principles that govern
 475 competition in free market economies will suffice to lead
 476 many institutions to adopt them.³ Schools and universities
 477 will make use of them in certain assessment contexts. Com-
 478 panies that believe that the involvement of human beings
 479 has a significant impact on the quality of their output will
 480 use them in new vetting procedures. Of course, AI content
 481 generators will continue to be *used* in all institutions: they
 482 provide a myriad of new productivity-enhancing methods.
 483 AI-generated content identifiers will simply be incorporated
 484 into institutions' existing methods for creating trust and pre-
 485 serving reputation. For instance, if a student submits work
 486 that is identified as AI-generated, the teacher may engage in
 487 additional interactions with the student, to check the content
 488 is understood; if a professional submits work identified as
 489 AI-generated, the assessor may likewise ask further ques-
 490 tions. The key idea is simply that AI-generated content must
 491 be treated in certain special ways, befitting its origin.

492 **Proposed rules for media companies**

493 As we also discussed in the first section, AI content genera-
 494 tion also allows people to *proliferate* content more than was
 495 previously possible, allowing content that is untethered from
 496 traditional human production processes to flow in large vol-
 497 umes into society. The mechanisms for disseminating con-
 498 tent in society can be thought of as the 'media', very broadly
 499 speaking, so we believe these organisations have important
 500 new roles in deploying reliable AI-generated content identi-
 501 fiers, if these are available. We will consider 'mainstream
 502 media' and 'social media' separately. We will also consider
 503 Web search companies, which are also involved in dissemin-
 504ating information.

505 **Mainstream media companies**

506 Mainstream media companies include traditional newspa-
 507 pers and radio and TV broadcasters. AI-generated content
 508 is finding its way into these venues in various forms: for
 509 instance in print articles (see e.g. Farhi, 2023), photos (see

³ We must of course ensure that identification methods are afford-
 able. We discuss the cost of identification methods later in the paper.

3FL01

3FL02

510 e.g. Oremus & Verma, 2023), and even video and audio
 511 content (see e.g. Stokel-Walker, 2023).

512 Mainstream media providers' business models certainly
 513 rely on reputation and trust, and we presume most such providers
 514 only include AI-generated content unintentionally. These providers
 515 certainly have an interest in using reliable
 516 AI-generated content identifiers if they are available. But
 517 many mainstream media providers are proving to be slow
 518 in adapting to the new AI world, and could benefit from
 519 guidance. Given that these providers disseminate content
 520 in large volumes to the wider public, we suggest they have
 521 a moral duty to use reliable content identifiers when these
 522 are available—and to use them systematically, so that *all*
 523 content they disseminate is checked. If content identifiers
 524 are affordable and run automatically, this filter should be
 525 minimally intrusive for companies—and would help to pre-
 526 serve their reputation in a world where AI-generated content
 527 is proliferating.

528 In most cases, we think it should be possible for media
 529 companies to disseminate AI-generated content, *if this is*
 530 *clearly flagged as such*. A flag would indicate, minimally,
 531 that the media outlet is *aware* that the flagged content is AI-
 532 generated, and can therefore be expected to have undertaken
 533 the kind of actions needed to preserve its reputation as a
 534 trustworthy provider. In fact there are some new companies
 535 that explicitly position themselves as providers of AI-generated
 536 content—in particular for local news: see for example
 537 NewsCorp's Data Local (Meade, 2023) and the UK's Radar
 538 News. The important thing is that these companies indicate
 539 clearly to their consumers that their content is AI-generated.
 540 The obligation to treat this content with due caution then
 541 falls on those who consume this content.

542 There may be some types of AI content where stronger
 543 obligations are appropriate. For instance, the Paris Charter
 544 on AI and Journalism (PAIJ, 2023) takes a stronger line
 545 on multimodal content 'mimicking real-world captures and
 546 recordings or realistically impersonating actual individuals'.
 547 The Charter recommends that outlets should *refrain* from
 548 using content of this kind. This proposed policy draws a very
 549 clear line between authentically captured content and syn-
 550 thetically created content. We feel that stronger moderation
 551 policies may indeed be required for AI content that convinc-
 552 ingly appears to have been recorded directly from the world.

553 If media providers have a moral duty to check for and
 554 appropriately moderate AI-generated content, we can ask
 555 whether this duty should also be encoded in law. It is likely
 556 that different jurisdictions will take different approaches
 557 here. For instance, US law places strong emphasis on free-
 558 dom of the press, while laws in European countries often
 559 define conditions on this freedom (see e.g. Tenorio, 2013).
 560 But the practical outcomes of press regulation are often
 561 more similar across jurisdictions than one might think
 562 (see e.g. Heller & van Hoboken, 2019): for instance, child

pornography is illegal everywhere. Clearly, the category of
 563 AI-generated content would require a much more nuanced
 564 moderation policy. Nonetheless, we believe there may be
 565 mechanisms in many jurisdictions for encoding rules about
 566 AI-generated content, and we recommend policymakers
 567 consider such rules.

568 In relation to existing rules: the EU's AI Act does in fact
 569 envisage a 'disclosure obligation' on the publishers of 'AI-
 570 generated or manipulated text' (in Recital 70b). This obliga-
 571 tion appears to be waived if the AI content 'has undergone
 572 a process of human review or editorial control and a natural
 573 or legal person holds editorial responsibility for the publica-
 574 tion of the content'. We think even in this case, there should
 575 be an obligation of some kind (whether legal or ethical) to
 576 explicitly flag AI-generated content. This is partly because
 577 'human review' is an imprecise concept: it's hard to know
 578 how engaged the human reviewer was in the process, espe-
 579 cially if large amounts of AI content are to be reviewed,
 580 because of the risk of 'automation bias' (see again Zerilli
 581 et al., 2019). But we also feel consumers have a right to
 582 know how much AI-generated content they are seeing: in
 583 other words, to know what the editorial practices on this
 584 matter are, for a given outlet.

Social media companies

586 Social media companies' business model is different from
 587 that of mainstream media companies. They both have incen-
 588 tives to maximise the viewer/user base; but social media
 589 companies have less incentive to present themselves as
 590 trusted information providers. Famously, under Section. 230
 591 of the US Communications Decency Act, social media com-
 592 panies are not responsible for the content they disseminate:
 593 rather, platform users have responsibility for the content they
 594 post. Individual users have incentives to disseminate AI-
 595 generated content, to increase the volume of content they
 596 produce. This could be motivated on financial grounds,
 597 to increase revenue from advertising, or simply through a
 598 desire to reach a large audience, to promote a political mes-
 599 sage, for instance. Reputation for individual users in this
 600 latter case is less of an issue, because users on social media
 601 are somewhat anonymous: it is easy for an individual to
 602 create multiple accounts, or to migrate between accounts,
 603 even if these practices are discouraged by most platforms.
 604 This means that large volumes of AI-generated content are
 605 likely to proliferate on social media platforms, as uptake of
 606 generators becomes a common public practice.

607 These considerations again lead us to recommend that
 608 social media companies should be *required* to use reliable
 609 AI-generated content identifiers when these are available,
 610 to systematically vet all content posted on their platforms,
 611 and moderate AI-generated content appropriately when
 612 it is found. We believe this is a crucial new regulatory

614 requirement, with an important role in preventing the dis-
 615 semination of content that is unconnected to traditional
 616 human production mechanisms, and an important role in
 617 extending society's existing mechanisms for regulating
 618 human communication into the new domain of AI-generated
 619 content.

620 **Web search companies**

621 Another important type of AI-content provider is 'fully AI-
 622 generated' websites. These are websites which are set up to
 623 cheaply disseminate information, in the interest of attracting
 624 users visiting from search engines (see e.g. Ryan-Mosley,
 625 2023). They exist independently on the Web, rather than
 626 within a social media platform. The relevant actors for *iden-*
 627 *tifying* AI-generated content in this case are Web search
 628 companies.

629 It is important that search engines deploy any reliable AI
 630 identification methods that exist, to systematically look for
 631 AI-generated sites, and inform their users of any sites that
 632 are found, whether by flagging identified sites or downrank-
 633 ing them in search results. We believe that the search engine
 634 companies are intrinsically motivated to do this, to retain
 635 the trust of their users. In this sense, the free market cre-
 636 ates incentives to use AI-content identifiers, as in the cases
 637 discussed above. But competition among search engines is
 638 not always strong; Google is still the dominant market leader
 639 (Oberlo, 2024). So we suggest policymakers should monitor
 640 whether free market considerations are sufficient to motivate
 641 search companies to make good use of AI content-identifi-
 642 cation resources. The EU's Digital Markets Act (EU, 2022)
 643 should enable this kind of monitoring, at least for search
 644 companies operating within the EU.

645 **How should media companies moderate 646 the AI-generated content they identify?**

647 Moderation methods are different for different types of
 648 media provider, so we will consider them separately. But we
 649 suggest one general rule for all providers: any content that
 650 is disseminated (or linked) that is identified as AI-generated
 651 should be clearly flagged as such.

652 **Mainstream media companies**

653 For mainstream media companies, the decision to publish
 654 a piece of AI-generated content will be taken by a human
 655 editor. Editors should certainly be able to run AI-generated
 656 content if they choose, as already noted. The key question
 657 is how to flag such content when it is published. There are
 658 various options to be explored. A textual flag could suf-
 659 fice, provided it is presented prominently enough to alert
 660 the consumer. A graphical flag could also be designed, that

conventionally denotes AI-generated content: perhaps an
 661 image of a robot with a pen.

663 **Social media companies**

664 For social media companies, decisions in relation to AI-
 665 generated content fall within the domain of content modera-
 666 tion. Content moderation methods on social media platforms
 667 involve many automated classifiers, looking for content of
 668 different kinds. Some moderation actions are taken auto-
 669 matically; others are passed to human moderators for final
 670 decisions. We recommend that AI-content detectors are
 671 incorporated into these moderation processes, to implement
 672 the following policy.

673 In the case where a single individual or group creates
 674 multiple accounts ('burner accounts'), that all disseminate
 675 AI-generated content pursuing a single goal, we recommend
 676 the appropriate moderation action is to remove this coordi-
 677 nated set of accounts altogether. This already seems to be
 678 standard policy for several social media platforms, such as
 679 Meta (see e.g. Facebook, 2023). Obviously the usual provi-
 680 sions for challenges and transparency should apply in such
 681 cases, as they do whenever an account is deleted.

682 In the case where a single user posts AI-generated con-
 683 tent, we suggest the content can always be left in place,
 684 provided it does not violate other company policies. But it
 685 should again be clearly flagged as AI-generated. For users
 686 who are posting large amounts of AI-generated content, for
 687 the sole purposes of increasing user engagement and adver-
 688 tising revenue, we suggest a further measure: content from
 689 such users should be downranked in platform recommender
 690 algorithms, so it disseminates less rapidly than other types
 691 of content. The amount of downranking of content from a
 692 given user could be a function of the amount of AI-generated
 693 content they are posting. (More generally, there could be
 694 limits imposed on the volume of AI-content disseminated by
 695 the platform as a whole, similar to the limits on the amount
 696 of pollution that can be produced by heavy industry.)

697 In addition to the above moderation policies (or perhaps
 698 instead of them), we suggest social media users should have
 699 broader agency of their own in relation to AI-generated con-
 700 tent. We suggest users should be able to configure settings
 701 for their own account so they can opt out of receiving *any*
 702 content that has been reliably identified as AI-generated,
 703 whatever its source. An alternative measure would be to
 704 allow users to *opt in* to receiving AI-generated content,
 705 so the default policy is that they receive none. The right
 706 choices here will depend on balancing the risks inherent
 707 in AI content moderation against those resulting from the
 708 unmoderated dissemination of AI content. We discuss how
 709 to approach this in the next subsection.

710 Finally, we suggest that social media companies have cer-
 711 tain new obligations in their reports to the general public,

712 if reliable AI content detection methods exist. They should
 713 report the overall amount of AI-generated content on their
 714 platforms, as part of regular transparency reporting. They
 715 should also report fluctuations in this amount, which may
 716 be linked to elections or other political events. And they
 717 should report the proportion of AI-generated content they
 718 removed—as well as the proportion of users who opted in
 719 (or out) of receiving AI-generated content, if these options
 720 are available. These reports are important in timely identifi-
 721 cation of risks arising from misinformation.

722 Web search companies

723 Web search companies already have mature policies that
 724 withhold or downrank content from untrusted providers.
 725 We suggest that AI-generated content should feature within
 726 these policies. In particular, websites that provide large
 727 amounts of AI-generated content, and do not clearly identify
 728 this content as AI generated, should be withheld from search
 729 results.⁴ Websites which occupy the ‘borderline’ on this cri-
 730 terion should be downranked in the search results. Google’s
 731 current stated policy is to rank content by quality, without
 732 regard for its human or AI origin (see e.g. Schwartz, 2024;
 733 Tucker, 2024). But there are likely already penalties for AI
 734 content that is presented deceptively as human-generated. If
 735 there aren’t, we suggest there should be.

736 In order to have some oversight over policies of this kind,
 737 as with social media companies, we also suggest that search
 738 companies should be required to report the overall amount
 739 of AI-generated content they identify on the Web, as part of
 740 their regular transparency reporting. Again, the EU’s Digital
 741 Markets Act may provide helpful mechanisms of overseeing
 742 this reporting.

743 Communication when AI-content detection is unreliable

744 In all the above policies, it is important to cater for circum-
 745 stances when reliable AI-content detection mechanisms are
 746 not available. In such contexts, the absence of an ‘AI-gener-
 747 ated’ flag on a piece of content does not positively indicate
 748 it is human-generated—and consumers need to know this.
 749 We suggest that in such situations, media companies display
 750 a general message for users, indicating that normal methods

751 for moderating AI-generated content are not running, or are
 752 impaired. This may be presented in some prominent place
 753 in a newspaper, or on the user’s app screen.

754 Balancing the risks of AI-content moderation 755 against the risks of AI-content proliferation

756 In any discussion of automated tools for identifying AI-gener-
 757 ated content, it is vital to consider the effects of *errors* in
 758 tool performance. We are aiming for ‘reliable’ tools, but in
 759 practice errors will always occur, and they can be harmful.
 760 False positives, where human-generated content is wrongly
 761 identified as AI-generated, are particularly harmful—at
 762 least, in that they create harms to the reputation of indi-
 763 vidual human generators of content, and may also infringe
 764 their rights to free expression, if identification triggers mod-
 765 eration actions. False negatives are also harmful, of course
 766 in misleading content consumers. How can these harms be
 767 balanced against the risks of unmoderated proliferation of
 768 AI-generated content? We suggest the main focus should
 769 be on minimising false positives. It will also be important
 770 to check for biases in false positives: we do not want to see
 771 more false positives for some demographic groups than oth-
 772 ers. There is clearly a need for discussion between agencies
 773 and providers as to what counts as a ‘reliable’ identification
 774 method. In relation to the EU’s AI Act, this will likely be
 775 decided as a technical standard, rather than in black-letter
 776 law, because the appropriate definition is likely to change
 777 as technologies advance.

778 Another important question concerns what stance to take
 779 for content that is generated partly by humans and partly
 780 by AI. For instance, if a user writes a text then asks GPT to
 781 ‘tidy it up’, we would not want this to be identified as a piece
 782 of ‘AI-generated content’. It is difficult to identify mixed
 783 human-LLM text using a classifier running externally to
 784 the provider company (see e.g. Gao et al., 2024). Detection
 785 methods that rely on company support have a strong advan-
 786 tage here, because they can make reference to the context in
 787 which the content was generated, including (crucially) the
 788 prompt history that led to the generated item. For instance,
 789 a company can choose to omit the identifying watermark or
 790 provenance metadata in cases where the human had a size-
 791 able role in creating the content—or to omit the generated
 792 content from the logged content, if a log-based detector is
 793 implemented.

794 A final important consideration in any discussion of con-
 795 tent moderation is freedom of speech. As a general rule,
 796 moderating content provided by a person infringes their
 797 right to freedom of expression if he/she does not give clear
 798 consent to the moderator. This is a fundamental human
 799 right—though of course, the right to freedom of expres-
 800 sion often trades off against other human rights (see e.g.
 801 Heyman, 1998). But in the case of AI-generated content,

4FL01⁴ A more far-reaching idea, which goes beyond the scope of the
 4FL02 current paper, is that a cap could be imposed on the amount of AI-
 4FL03 generated content a single provider can make available. The idea of
 4FL04 capping ‘volume’ of content has precedents in other areas of regula-
 4FL05 tion—for instance, in the regulation of polluters. A rule of this kind
 4FL06 may be useful in addressing wider problems of information overload
 4FL07 (see e.g. Holyst et al., 2024). Such a rule could potentially make use
 4FL08 of an AI content detection tool—but it might more practically be
 4FL09 enforced by restrictions on compute resources allocated to companies
 4FL10 (see Sastry et al., 2024 for a relevant proposal).

802 some completely new considerations may arise. If Joe posts
 803 a piece of content that was produced (from scratch) *by an AI*
 804 *system*, and this content is moderated, is Joe's right to free
 805 expression in any way being curtailed? Ex hypothesis, Joe
 806 did not *express* the content. Joe *disseminated* it (by posting
 807 it), but he didn't create it. Of course, there are gradations of
 808 human involvement in AI content generation, as just dis-
 809 cussed: the more involved Joe is in the process, the more
 810 rights he has. The act of posting content can likewise involve
 811 gradations of human involvement. Nonetheless, the concept
 812 of freedom of expression may apply somewhat differently
 813 to AI-generated content—arguably removing some of the
 814 difficult issues that arise in most content moderation. The
 815 strong moderation actions we recommended above for media
 816 companies all apply in cases where the human provider is
 817 minimally involved, or not involved at all, and particularly
 818 if the provider is anonymous.

819 **Support for reliable identification 820 mechanisms in the wider tech world**

821 In the previous section, we asked how reliable methods for
 822 identifying AI-generated content should be deployed, if they
 823 are available. But as discussed in the section before that, we
 824 find ourselves in a new adversarial situation, in which some
 825 actors have incentives to defeat the dominant identification
 826 methods. In this section, we conclude by considering what
 827 policies would help give identification methods the upper
 828 hand in this new arms race. Of course, we can learn a lot
 829 from long-running arms races in other areas—for instance,
 830 relating to search engine optimisation or malicious content
 831 detection. In particular, techniques for identifying coordi-
 832 nated malicious efforts (see e.g. Pacheco et al., 2021) will
 833 readily extend to AI-fuelled disinformation campaigns. But
 834 the AI-content-detection arms race also offers new technical
 835 opportunities for interventions, because the adversarial
 836 content in this case is all AI-generated. In this section, we
 837 review these new opportunities.

838 **Regulation on provenance-authentication protocols**

839 As we noted earlier, requiring the providers of AI content
 840 generators to support detection only covers *one* method
 841 for identifying AI-generated content. Another method
 842 involves establishing broader protocols for provenance
 843 authentication, that apply to human-generated content as
 844 well as AI-generated content. Through these protocols,
 845 trusted providers of content, whether AI-generated or
 846 human-generated, can positively identify the content they
 847 provide. Content whose provenance is *not* authenticated
 848 can then be regarded with more caution, and perhaps

849 moderated accordingly. The details of a workable prove-
 850 nance-authentication scheme still remain to be worked out:
 851 implementing such a scheme is a long term project. In par-
 852 ticular, it is important to implement a way of authenticat-
 853 ing content as produced by an individual person, without
 854 disclosing this person's identity. (A system such as that
 855 used for German ID cards is one possibility here; see e.g.
 856 Poller et al., 2012.)

857 We also noted earlier that provenance authentication
 858 mechanisms require support throughout the information
 859 ecosystem, from creation and capture, through transmis-
 860 sion and modification, to final display. So if there is to be
 861 regulation in this area, it must be separate from regula-
 862 tion focussed narrowly on AI providers. In this section,
 863 we will consider possible regulatory actions relating to
 864 provenance-authentication.

865 Our main point is that rules requiring AI providers to
 866 support content detection and rules requiring the wider
 867 ecosystem to adopt provenance methods should not be
 868 seen as alternatives to one another. We see roles for both
 869 types of rule. Crucially, neither type of rule provides a
 870 failsafe method for the identification of AI-generated con-
 871 tent, in the arms race we are embarking on. As we already
 872 stressed above, the rules in the AI Act will sometimes be
 873 defeated by adversaries, will be flatly ignored by mali-
 874 cious actors, and will not thoroughly permeate the open-
 875 source generator ecosystem. A provenance scheme pro-
 876 vides a good supplement to detector tools. Conversely, a
 877 provenance-authentication scheme is also fallible, and has
 878 important limits. For instance, authentication information
 879 can often be removed or changed if a piece of content is
 880 copied. It will also be difficult to instrument every device
 881 that can manipulate content.

882 As already noted, voluntary schemes for adopting
 883 provenance protocols are already beginning to infiltrate
 884 the tech world. But widespread adoption is necessary to
 885 ensure the success of a provenance scheme. We believe
 886 this will only be possible if broader legislation supporting
 887 provenance-authentication is enacted. But crucially, this
 888 broader legislation should complement legislation requir-
 889 ing providers of AI content generators to support detection
 890 mechanisms.

891 Once again, the EU's AI Act is very well formulated to
 892 accommodate provenance authentication schemes. Recital
 893 70a, which states the context for rules on content identi-
 894 fication, makes reference to provenance schemes as well
 895 as to detection methods. But Article 52.1(a), which states
 896 the obligations on AI providers, refers only to support for
 897 detection methods. The Act would therefore dovetail well
 898 with additional broader rules about provenance authentica-
 899 tion. Biden's Executive Order also envisages a division of
 900 labour between detection schemes and provenance authen-
 901 tication schemes.

902 **Regulation preventing the open-sourcing
903 of 'frontier' AI models**

904 Enforcing regulations on AI systems is harder in the open-
905 source world than for proprietary commercial systems. For
906 instance, as we discussed earlier in the paper, the rule that
907 AI providers must support detection mechanisms is harder
908 to enforce for open-source AI generators than for com-
909 mercial generators. Copies of open-source generators can
910 proliferate, existing code supporting detection can be mod-
911 ified or removed. Open-source generators are also helpful
912 to actors looking for ways to evade detectors elsewhere
913 in the ecosystem: they provide a platform for exploring
914 evasion methods.

915 A debate is emerging between groups seeking to promote
916 the practice of open-sourcing generative AI models (such
917 as the AI Alliance) and groups seeking to prevent the prac-
918 tice: see Bommasani et al. (2023) for a good overview. In
919 relation to detection of AI-generated content, we see con-
920 siderable risks in the practice of open-sourcing generative
921 AI models—especially for the 'frontier' models with the
922 best performance, created by the best-resourced providers.
923 In this sense, we align ourselves with the recent stance of
924 Seger et al. (2023), who argue persuasively that many risks
925 arise from the open-sourcing of these frontier models. We
926 suggest that regulation that prevents the open-sourcing of
927 new frontier models (or in Seger's terms, 'highly capable'
928 AI models) will do a great deal to stack the playing field in
929 favour of reliable AI-content identification mechanisms. (A
930 recent analysis by Kapoor et al., 2024 also summarises risks
931 of open-source foundation models, but is more equivocal in
932 its conclusions.)

933 **Support for applied research in detection
934 mechanisms**

935 In the adversarial climate we sketched above, new or
936 extended detection mechanisms for AI-generated content
937 will always be needed. This research could come from
938 academia or from industry: in either case, there is a good
939 argument that governments should actively support such
940 research. Results from this research should perhaps be kept
941 out of public venues, if this would make it harder for new
942 schemes to be attacked.

943 **Support for compliance with identification schemes**

944 Rules requiring provenance-authentication schemes and
945 rules requiring AI providers to support detection schemes
946 obviously need to be enforced, in jurisdictions where they
947 apply. In these contexts, policymakers also have a role in

948 resourcing compliance and enforcement efforts, and making
949 enforcement as efficient as possible.

950 As regards compliance, it is vitally important to consider
951 the financial costs of complying with mandated detection or
952 provenance-authentication schemes—especially given the
953 importance of making identification methods available at
954 low costs (which we have already emphasised). We might
955 imagine governments bearing some of these costs—espe-
956 cially for smaller companies, for whom they would be par-
957 ticularly burdensome. At a national level, institutions like the
958 UK's new AI Safety Institute may have a role to play here.
959 International bodies could also have a role; for instance, the
960 EU's newly formed AI Office.

961 As regards efficiency, there are two useful directions.
962 Firstly, large providers of AI generators which are not pro-
963 viding all possible support for detection tools should be a
964 focus for enforcement. Part of the effort should be to dis-
965 seminate good information about the best available tools to
966 providers. Providers in the open-source community may be a
967 particular focus here. Secondly, certain links in the informa-
968 tion ecosystem have particular roles in attacks on AI-content
969 detection methods. For instance, as we have already dis-
970 cussed, systems that paraphrase text or alter images can be
971 used to evade detection. It is particularly important that these
972 content-modification systems adopt provenance protocols, to
973 provide relevant information to content consumers.

974 **Summary**

975 In this paper, we have sketched the problems that are likely
976 to arise if AI-generated content disseminates into society on
977 a large scale without appropriate checks and balances. We
978 have summarised some recent policy initiatives in the EU
979 and US that address this scenario, by requiring AI provid-
980 ers to support mechanisms that allow reliable identification
981 of AI-generated content. We applaud these new initiatives.
982 They are not a panacea, but we judge that they will apply a
983 consistent impetus on AI providers, to create reliable detec-
984 tion mechanisms. They create a new dynamic context, in
985 which policymakers can consider some new questions.

986 Our paper considers what new options there are for poli-
987 cymakers in this new dynamic context. Our recomme-
988 dations are of two types. Firstly, we recommend some new
989 rules about who should *use* reliable AI-content detectors,
990 when these are available, and how they should be used. Our
991 proposals here focus on new obligations for media compa-
992 nies. We make different recommendations for mainstream
993 media companies, social media companies and Web search
994 companies. Secondly, we recommend some new rules that
995 will help create an environment where reliable AI-generated
996 content identification methods exist. We suggest a vari-
997 ety of different rules: rules instituting broad protocols for

998 provenance-authentication throughout the digital information
 999 ecosystem; rules preventing the open-sourcing of new
 1000 'frontier' generative AI models; policies supporting applied
 1001 research in AI-generated content detection; and policies sup-
 1002 porting compliance with identification schemes, including
 1003 through assistance with costs of compliance.
 1004

1005 **Data availability** No datasets were generated or analysed in the study
 1006 reported in this paper.

1007 References

1008 Berry, S. (2024). Fake Google restaurant reviews and the implications
 1009 for consumers and restaurants. PhD dissertation, William Howard
 1010 Taft University. <https://arxiv.org/pdf/2401.11345.pdf>

1011 Bradford, A. (2020). *The Brussels effect: How the European Union
 1012 rules the world*. Oxford University Press.

1013 Candelier, F., Krayer, L., Rajendran, S. and Zuluaga Martínez, D.
 1014 (2023). How People Can Create—and Destroy—Value with
 1015 Generative AI. BCG Henderson Institute report. [https://www.bcg.com/publications/2023/how-people-create-and-destro-1017 oy-value-with-gen-ai](https://www.bcg.com/publications/2023/how-people-create-and-destro-1016 oy-value-with-gen-ai)

1017 Crothers, E., Japkowicz, N., & Viktor, H. L. (2023). Machine-gener-
 1018 ated text: A comprehensive survey of threat models and detection
 1019 methods. *IEEE Access*, 11, 70977–71002.

1020 Davis, H. (2006). Search engine optimization.

1021 Dell'Acqua, F., McFowland, E., Mollick, E. R., Lifshitz-Assaf, H.,
 1022 Kellogg, K., Rajendran, S., Krayer, L. Candelier, F., & Lakhani,
 1023 K. R. (2023). Navigating the jagged technological frontier: Field
 1024 experimental evidence of the effects of AI on knowledge worker
 1025 productivity and quality. Harvard Business School Technology &
 1026 Operations Mgt. Unit Working Paper, (24–013).

1027 de Wilde, P., Arora, P., Buarque de Lima Neto, F., Chin, Y., Thinyane,
 1028 M., Stinckwich, S., Fournier-Tombs, E., & Marwala, T. (2024).
 1029 *Recommendations on the use of synthetic data to trainAI models*.
 1030 United Nations University Policy Guideline. [https://collections.unu.edu/eserv/UNU:9480/Use-of-Synthetic-Data-to-Train-AI-1032 Models.pdf](https://collections.unu.edu/eserv/UNU:9480/Use-of-Synthetic-Data-to-Train-AI-1031 Models.pdf)

1033 Dohmatob, E., Feng, Y., & Kempe, J. (2024a). Model Collapse Demys-
 1034 tified: The Case of Regression. arXiv preprint [arXiv:2402.07712](https://arxiv.org/abs/2402.07712).

1035 Dohmatob, E., Feng, Y., Yang, P., Charton, F., & Kempe, J. (2024b). A
 1036 Tale of Tails: Model Collapse as a Change of Scaling Laws. arXiv
 1037 preprint [arXiv:2402.07043](https://arxiv.org/abs/2402.07043).

1038 EU (2022). Regulation (EU) 2022/1925 of the European Parliament
 1039 and of the Council of 14 September 2022 on contestable and
 1040 fair markets in the digital sector and amending Directives (EU)
 1041 2019/1937 and (EU) 2020/1828 (Digital Markets Act)". EUR-Lex.

1042 EU/FLI (2024). EU Artificial Intelligence Act. The Act Texts.
 1043 Resources provided by the Future of Life Institute. [https://artif-1045 icialintelligenceact.eu/the-act/](https://artif-1044 icialintelligenceact.eu/the-act/)

1046 Facebook (2023). Account integrity and authentic identity. Facebook
 1047 Transparency Center. [https://transparency.fb.com/en-gb/policies/1049 community-standards/account-integrity-and-authentic-identity/](https://transparency.fb.com/en-gb/policies/1048 community-standards/account-integrity-and-authentic-identity/)

1050 Farhi, P. (2023). A news site used AI to write articles. It was a jour-
 1051 nalistical disaster. Washington Post, January 2023. [https://www.wash-1054 ingtonpost.com/media/2023/01/17/cnet-ai-articles-journ-1055 alism-corrections/](https://www.wash-1052 ingtonpost.com/media/2023/01/17/cnet-ai-articles-journ-1053 alism-corrections/)

Fernandes, F (2023). Mapped: Interest in Generative AI by Country.
 Visual Capitalist blog post. [https://www.visualcapitalist.com/cp-1054 mapped-interest-in-generative-ai-by-country/](https://www.visualcapitalist.com/cp-1053 mapped-interest-in-generative-ai-by-country/)

FID (2024). AI as a Public Good: Ensuring Democratic Control of
 1055 AI in the Information Space. Report by the Forum for Infor-
 1056 mation and Democracy. [https://informationdemocracy.org/2024/02/1059 28/new-report-of-the-forum-more-than-200-policy-recom-1060 mendations-to-ensure-democratic-control-of-ai/](https://informationdemocracy.org/2024/02/1057 28/new-report-of-the-forum-more-than-200-policy-recom-1058 mendations-to-ensure-democratic-control-of-ai/)

Founta, A. M., Chatzakou, D., Kourtellis, N., Blackburn, J., Vakali, A.,
 1061 & Leontiadis, I. (2019, June). A unified deep learning architecture
 1062 for abuse detection. In Proceedings of the 10th ACM conference
 1063 on web science (pp. 105–114).

Gao, C., Chen, D., Zhang, Q., Huang, Y., Wan, Y., & Sun, L. (2024).
 1064 LLM-as-a-coauthor: The challenges of detecting LLM-human
 1065 mixcase. arXiv preprint [arXiv:2401.05952](https://arxiv.org/abs/2401.05952).

1066 Google (2024). Gemini 1.5: Unlocking multimodal understanding
 1067 across millions of tokens of context. arXiv preprint [arXiv:2403.05530](https://arxiv.org/abs/2403.05530).

HackerNoon (2023). AI Design Tools That are Changing How Graphic
 1068 Designers Work. [https://hackernoon.com/ai-design-tools-that-are-1070 derailing-how-graphic-designers-work](https://hackernoon.com/ai-design-tools-that-are-1069 derailing-how-graphic-designers-work)

Hans, A., Schwarzschild, A., Cherepanova, V., Kazemi, H., Saha,
 1071 A., Goldblum, M., Geiping, J., & Goldstein, T. (2024). Spotting
 1072 LLMs With Binoculars: Zero-Shot Detection of Machine-Gen-
 1073 erated Text. arXiv preprint [arXiv:2401.12070](https://arxiv.org/abs/2401.12070).

Heller, B., & van Hoboken, J. (2019). Freedom of expression: A com-
 1074 parative summary of United States and European law. Available
 1075 at SSRN 4563882. <https://doi.org/10.2139/ssrn.4563882>

1076 Heyman, S. J. (1998). Righting the balance: An inquiry into the foun-
 1077 dations and limits of freedom of expression. *BUL Rev*, 78, 1275.

Holyst, J. A., Mayr, P., Thelwall, M., Frommholz, I., Havlin, S., Sela,
 1078 A., & Sienkiewicz, J. (2024). Protect our environment from infor-
 1079 mation overload. *Nature Human Behaviour*, 8, 402–403.

Jakesch, Maurice, Advait Bhat, Daniel Buschek, Lior Zalmanson,
 1080 and Mor Naaman. 2023. Co-Writing with Opinionated Language
 1081 Models Affects Users' Views. In Proceedings of the 2023 CHI
 1082 Conference on Human Factors in Computing Systems, 1–15.
 1083 Hamburg, Germany: ACM.

Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bam-
 1084 ford, C., & Sayed, W. E. (2024). Mixtral of experts. arXiv preprint
 1085 [arXiv:2401.04088](https://arxiv.org/abs/2401.04088). <https://doi.org/10.48550/arXiv.2401.04088>

Kapoor, S., Bommastan, R., Klyman, K., Longpre, S., Ramaswami,
 1086 A., Cihon, P., Hopkins, A., Bankston, K., Biderman, S., Bogen,
 1087 M., Chowdhury, R., Engler, A., Henderson, P., Jernite, Y., Lazar,
 1088 S., Maffulli, S., Nelson, A., Pineau, J., Skowron, A., Song, D.,
 1089 Storchan, V., Zhang, D., Ho, D., Liang, P., Narayanan, A. (2024).
 1090 On the Societal Impact of Open Foundation Models. Stanford
 1091 University Center for Research on Foundation Models. <https://1092 crfm.stanford.edu/open-fms/paper.pdf>

Knott, A., Pedreschi, D., Chatila, R., Chakraborti, T., Leavy, S., Baeza-
 1093 Yates, R., Evers, D., Trotman, A., Teal, P. D., Biecek, P., Russell,
 1094 S., & Bengio, Y. (2023). Generative AI models should include
 1095 detection mechanisms as a condition for public release. *Ethics
 1096 and Information Technology*, 25(4), 55.

Krishna, K., Song, Y., Karpinska, M., Wieting, J., & Iyyer, M. (2023).
 1097 Paraphrasing evades detectors of AI-generated text, but retrieval is
 1098 an effective defense. *Advances in Neural Information Processing
 1099 Systems*. <https://arxiv.org/abs/2303.13408>

Liang, W., Izzo, Z., Zhang, Y., Lepp, H., Cao, H., Zhao, X., ... & Zou,
 1100 J. Y. (2024). Monitoring AI-Modified Content at Scale: A Case
 1101 Study on the Impact of ChatGPT on AI Conference Peer Reviews.
 1102 arXiv preprint [arXiv:2403.07183](https://arxiv.org/abs/2403.07183).

Májovský, M., Černý, M., Netuka, D., & Mikolov, T. (2024). Perfect
 1103 detection of computer-generated text faces fundamental chal-
 1104 lenges. *Cell Reports Physical Science*, 5(1), 101769.

Meade, C. (2023). News Corp using AI to produce 3,000 Australian
 1105 local news stories a week. The Guardian, July 2023. [https://www.theguardian.com/media/2023/aug/01/news-corp-ai-chat-gpt-sto-1107 ries](https://www.theguardian.com/media/2023/aug/01/news-corp-ai-chat-gpt-sto-1106 ries)

1108

1122 Munich (2024). Tech Accord to Combat Deceptive Use of AI in 2024
1123 Elections. Pledge made at the Munich Security Conference, February 2024. <https://securityconference.org/en/ai-elections-accord/>

1124 NBC (2024). Fake Joe Biden robocall tells New Hampshire Democrats
1125 not to vote Tuesday. NBC News. <https://www.nbcnews.com/politics/2024-election/fake-joe-biden-robocall-tells-new-hampshire-democrats-not-vote-tuesday-rcna134984>.

1126 Newsguard (2024). Tracking AI-enabled Misinformation: 702 'Unreliable
1127 AI-Generated News' Websites (and Counting). <https://www.newsguardtech.com/special-reports/ai-tracking-center/>

1128 Notopoulos, K. (2024). Women laughing alone with AI-generated content
1129 spam. Business Insider <https://www.businessinsider.com/the-hairpin-blog-ai-spam-content-farm-cybersquatting-2024-1>

1130 Oberlo (2024). Search Engine Market Share in 2024. <https://www.oberlo.com/statistics/search-engine-market-share>

1131 OpenAI (2021). DALL-E: creating images from text. Retrieved from
1132 <https://openai.com/research/dall-e> (accessed 19 March 2024).

1133 OpenAI. (2023). GPT-4: Scaling up deep learning. Retrieved from
1134 <https://openai.com/research/gpt-4>

1135 OpenAI. (2024). Sora: Creating video from text. Retrieved from <https://openai.com/sora>

1136 Oremus, W and Verma, P. These look like prizewinning photos.
1137 They're AI fakes. Washington Post, November 2023. <https://www.washingtonpost.com/technology/2023/11/23/stock-photos-ai-images-controversy/>

1138 Originality (2024). Ai-generated Research Papers Published On Arxiv
1139 Post Chatgpt Launch. Originality.AI blog post. <https://originality.ai/blog/ai-generated-research-papers>

1140 Pacheco, D., Hui, P.-M., Torres-Lugo, C., Truong, B. T., Flammini, A.,
1141 & Menczer, F. (2021). Uncovering coordinated networks on social
1142 media: Methods and case studies. *Proceedings of the International
1143 AAAI Conference on Web and Social Media*, 15(1), 455–466.

1144 PAIJ (2023). Paris Charter on AI and Journalism. <https://rsf.org/sites/default/files/medias/file/2023/11/Paris%20Charter%20on%20AI%20and%20Journalism.pdf>

1145 Poller, A., Waldmann, U., Vowé, S., & Türpe, S. (2012). Electronic
1146 identity cards for user authentication—promise and practice. *IEEE
1147 Security & Privacy Magazine*, 10(1), 46–54.

1148 Rawte, V., Sheth, A., & Das, A. (2023). A survey of hallucination
1149 in large foundation models. *arXiv preprint arXiv: 2309.05922*.
<https://doi.org/10.48550/arXiv.2309.05922>

1150 Rishi Bommasani, Sayash Kapoor, Kevin Klyman, Shayne Longpre,
1151 Ashwin Ramaswami, Daniel Zhang, Marietje Schaake, Daniel E.
1152 Ho, Arvind Narayanan, Percy Liang (2023). Considerations for
1153 Governing Open Foundation Models. Stanford University Center
1154 for Research on Foundation Models.

1155 Ryan-Mosley, T. (2023). Junk websites filled with AI-generated text
1156 are pulling in money from programmatic ads. MIT Technology
1157 Review. <https://www.technologyreview.com/2023/06/26/1075504/junk-websites-filled-with-ai-generated-text-are-pulling-in-money-from-programmatic-ads/>

1158 Sadasivan, V. S., Kumar, A., Balasubramanian, S., Wang, W., & Feizi,
1159 S. (2023). Can AI-generated text be reliably detected? *arXiv preprint arXiv: 2303.11156*. <https://doi.org/10.48550/arXiv.2303.11156>

1160 Sastry, G., Heim, L., Belfield, H., Anderljung, M., Brundage, M.,
1161 Hazell, J., & Coyle, D. (2024). Computing power and the governance
1162 of artificial intelligence. *arXiv preprint arXiv: 2402.08797*.
<https://doi.org/10.48550/arXiv.2402.08797>

1163 Schwartz, B. (2024). Google Responds To Claims Of Google News
1164 Boosting Garbage AI Content. Search Engine Roundtable, Jan
1165 2024. <https://www.seroundtable.com/google-responds-garbage-ai-content-in-google-news-36757.html>

1166 Seger, E., Dreksler, N., Moulange, R., Dardaman, E., Schuett, J., Wei,
1167 K., Gupta, A. (2023). Open-Sourcing Highly Capable Foundation
1168 Models: An Evaluation of Risks, Benefits, and Alternative
1169 Methods for Pursuing Open-Source Objectives.

1170 Shumailov, I., Shumaylov, Z., Zhao, Y., Gal, Y., Papernot, N., & Anderson, R. (2023). The curse of recursion: Training on generated data
1171 makes models forget. *arXiv preprint arXiv: 2305.17493*. <https://doi.org/10.48550/arXiv.2305.17493>

1172 Srinivasan, S. (2024). Detecting AI fingerprints: A guide to watermarking
1173 and beyond. Brookings Institute report. <https://www.brookings.edu/articles/detecting-ai-fingerprints-a-guide-to-watermarking-and-beyond/>

1174 Stokel-Walker, C. (2023). TV channels are using AI-generated presenters
1175 to read the news. The question is, will we trust them? BBC
1176 News, January 2024. <https://www.bbc.com/future/article/20240126-ai-news-anchors-why-audiences-might-find-digitally-generated-tv-presenters-hard-to-trust>

1177 Su, J., Zhuo, T. Y., Wang, D., & Nakov, P. (2023). DetectLLM: Leveraging
1178 Log Rank Information for Zero-Shot Detection of Machine-
1179 Generated Text. *arXiv preprint arXiv:2306.05540*.

1180 Tenorio, P. (2013). Freedom of Communication in the US and Europe.
1181 *ICL Journal*, 7(2), 150–173.

1182 UAE TII. Falcon-180b: A 180 billion token language model. <https://huggingface.co/tiuae/falcon-180B>, 2023.

1183 Tonmoy, S. M., Zaman, S. M., Jain, V., Rani, A., Rawte, V., Chadha,
1184 A., & Das, A. (2024). A comprehensive survey of hallucination
1185 mitigation techniques in large language models. *arXiv preprint arXiv:2401.01313*.

1186 Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M. A.,
1187 Lacroix, T., ... & Lample, G. (2023). Llama: Open and efficient
1188 foundation language models. *arXiv preprint arXiv:2302.13971*.

1189 Tucker, E. (2024). New ways we're tackling spammy, low-quality
1190 content on Search. Google blog post, March 2024. <https://blog.google/products/search/google-search-update-march-2024/>

1191 Valyaeva, I (2023). AI Has Already Created As Many Images As Photographers Have Taken in 150 Years. Statistics for 2023. EveryPixel Journal. <https://journal.everypixel.com/ai-image-statistics>

1192 Veselovsky, V., Ribeiro, M. H., & West, R. (2023). Artificial Artificial
1193 Artificial Intelligence: Crowd Workers Widely Use Large Language
1194 Models for Text Production Tasks. *arXiv preprint arXiv: 2306.07899*.

1195 Wang, Z., Bao, J., Zhou, W., Wang, W., Hezhen, Hu., Chen, H., & Li,
1196 H. (2023). DIRE for diffusion-generated image detection. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023, 22445–22455.

1197 Zerilli, J., Knott, A., Maclaurin, J., & Gavaghan, C. (2019). Algorithmic DECISION-MAKING AND THE CONTROL PROBLEM.
1198 *Minds and Machines*, 29, 555–578.

1199 Yichi Zhang and Xiaogang Xu. Diffusion noise feature: Accurate and
1200 fast generated image detection. *arXiv preprint arXiv:2312.02625*,
1201 2023.

1202 Zhou, Z. H. (2014). Ensemble methods. *Combining pattern classifiers*
1203 (pp. 186–229). Wiley.

1204 **Publisher's Note** Springer Nature remains neutral with regard to
1205 jurisdictional claims in published maps and institutional affiliations.

1206 Springer Nature or its licensor (e.g. a society or other partner) holds
1207 exclusive rights to this article under a publishing agreement with the
1208 author(s) or other rightsholder(s); author self-archiving of the accepted
1209 manuscript version of this article is solely governed by the terms of
1210 such publishing agreement and applicable law.