



**CONUL Regulatory Affairs
Research Paper**

**Artificial Intelligence and digital content:
the emerging impact of policy, law and practice
on higher education and research libraries**

14 April 2025

Table of Contents

1. Introduction	3
1.1. A focus on generative AI	3
1.2. Content in AI	4
1.3. AI tools	6
2. Ethical aspects	7
3. Policy and regulation	9
3.1. The European Union - a legislative approach	10
3.2. The United States - litigating change?	14
3.3. The United Kingdom - legal evolution	18
3.4. General legal directions of travel	20
4. Implications for higher education and research libraries	22
4.1. The higher education sector	22
4.2. AI products in education and research	25
4.3. AI and research libraries	26
4.4. AI and publishing of research	27
5. Conclusions	28

Acknowledgements

CONUL acknowledges the contributions of the following members of Regulatory Affairs in the production of this research paper:

David Meehan (DCU) (principal author); and

Aisling Keane (UG); Della Keating (NLI); Gillian Kerins (TUD); Elizabeth Murphy (MU); Emer Twomey (UCC); Kathryn Smith (RCSI); Michelle Dalton (UCD).

Copyright statement

© Consortium of National University Libraries (CONUL), 2025

This work is licensed under the Creative Commons Attribution Non-commercial No-derivatives 4.0 International licence (CC BY-NC-ND 4.0). For details go to <https://creativecommons.org/licenses/bync-nd/4.0/deed.en>

Please attribute using the following: Consortium of National and University Libraries (CONUL), 2025. Research Paper. Artificial Intelligence and digital content: the emerging impact of policy and law on higher education and research libraries. David Meehan and CONUL Regulatory Affairs. Ireland.

Disclaimer

The contents of this paper do not comprise or represent legal advice by CONUL or the paper's author or contributors. For formal legal guidance, please refer to the person or entity in your organisation responsible for obtaining, or referring for, legal services.

1. Introduction

Artificial intelligence (AI) is having a profound impact on the manner in which information is being generated, marketed and exploited. It is potentially the most disruptive phenomenon to emerge from the information technology revolution, encroaching on established institutions operating in business, scientific discovery, publishing, media and education, and altering the research environment for innovators, authors, academics and students. It promises advances in efficiency, productivity and convenience, while raising concerns for employment and the conduct of business, and, in the most extreme scenarios, consequences for humanity itself.

This paper will give some direction to stakeholders in higher education (HE) and research institutions and their libraries on ethical, policy and regulatory aspects of AI, and outline some impacts on the wider research environment. We are assuming that AI tools are already routinely deployed by staff, researchers and students in HE, and that its use will continue to proliferate.

1.1. A focus on generative AI

We will concentrate primarily on generative AI, and specifically its text mode, as opposed to images, audio, video and other content forms, although developments concerning the latter will be referred to where relevant. Generative AI analyses and connects data with its system to produce new data in response to user 'prompts'. These systems, or models, can be designed to perform simple or complex roles, with the larger models requiring greater investment in data acquisition and processing power.

Generative AI can be thought of as a machine-learning model trained to create new data, rather than making a prediction about a specific dataset. It draws on research and computational advances going back over 50 years.¹ It has been described as 'weak' or 'narrow' AI 'trained to perform a single or narrow task, often far faster and better than a human mind can'.

Higher capability artificial general intelligence (AGI), also known as strong AI, is understood as having the ability to reason and to exceed human cognitive capability. Described on the one hand as 'nothing more than a theoretical concept' and, even, dumber than a house cat,² and on the other as having the potential to harm

¹ Adam Zewe. 'Explained: Generative AI: How do powerful generative AI systems like ChatGPT work, and what makes them different from other types of artificial intelligence?'. *MIT News*, 9 November 2023. See: <https://news.mit.edu/2023/explained-generative-ai-1109> (accessed 14 April 2025).

² IBM, 'Understanding the different types of artificial intelligence', 12 October 2023, See: <https://www.ibm.com/think/topics/artificial-intelligence-types> (accessed 14 April 2025). Christopher Mims. 'This AI pioneer thinks AI is dumber than a cat'. *Wall Street Journal*, 11 October 2024.

humanity,³ it has been claimed that the cognitive threshold has already been approached, if not breached, by OpenAI's GPT-4.⁴

For many looking at short to mid-term trajectories, AI is being marketed as a tool that will not of itself replace people, but as one where individuals using it will replace those who don't.⁵ While some worry about the rate of progress towards generative AI,⁶ others note that while GPT-4 was a step-change over GPT-3.5, projected advances to next-stage GPT-5 (especially OpenAI's Orion, which had been expected to be feasible by mid-2024), have stalled, with the costs of training exceeding the benefits, and questions as to whether there is enough data available to support an upgraded model, or even whether continual scaling up of training data and infrastructure is beneficial in the first place.⁷ However, although they won't be considered further here, claims are being made on social media of progress from predictive (generative) AI to functioning reflective or reasoning AGI models, such as through OpenAI's o1 and o3 models released in December 2024 and January 2025 respectively.⁸

1.2. Content in AI

The current basis of AI is the use of 'Large Language Model' (LLM) technology to acquire and process content with a view to generating output flowing from 'prompts'

³ A reported claim by Geoffrey Hinton. See: Alyssa Lukpat. 'AI employees fear they aren't free to voice their concerns'. *Wall Street Journal*, 4 June 2024; and Aylin Woodward. 'Geoffrey Hinton, godfather of AI who expressed alarm over the technology, shares Nobel Prize in physics'. *Wall Street Journal*, 8 October 2024.

⁴ Lauren Leffer. 'In the Race to Artificial General Intelligence, Where's the Finish Line?'. *Scientific American*, 25 June 2024. See: <https://www.scientificamerican.com/article/what-does-artificial-general-intelligence-actually-mean/> (accessed 14 April 2025).

⁵ Karim Lakhani. 'AI Won't Replace Humans - But Humans With AI Will Replace Humans Without AI'. *Harvard Business Review*, 4 August 2023. See: <https://hbr.org/2023/08/ai-wont-replace-humans-but-humans-with-ai-will-replace-humans-without-ai> (accessed 14 April 2025).

⁶ Zoe Kleinman, Chris Vallance. 'AI 'godfather' Geoffrey Hinton warns of dangers as he quits Google'. *BBC News*, 2 May 2023. See: <https://www.bbc.com/news/world-us-canada-65452940> (accessed 14 April 2025).

⁷ Deepa Seetharaman. 'The Next Great Leap in AI Is Behind Schedule and Crazy Expensive'. *Wall Street Journal*, 20 December 2024; and Jeremy Hsu. 'AI scientists are sceptical that modern models will lead to AGI'. *New Scientist*, 14 March 2025, available at: <https://www.newscientist.com/article/2471759-ai-scientists-are-sceptical-that-modern-models-will-lead-to-agi/> (accessed 14 April 2025). A leveling off of gains in training to achieve effective reasoning was noted in a study using Claude 3.7. See: 'Reasoning models don't always say what they think'. 3 April 2025 at <https://www.anthropic.com/research/reasoning-models-dont-say-think> and a full paper at: https://assets.anthropic.com/m/71876fabef0f0ed4/original/reasoning_models_paper.pdf (both accessed 14 April 2025).

⁸ See <https://openai.com/o1/> (accessed 14 April 2025), and a general podcast discussion at: Dónal Mulligan, Ciarán O'Connor. 'Political upheaval and reasoning models'. *Enough about AI*, Episode 1, 2025, 24 February 2025.

devised by users. Strictly speaking, a properly administered LLM is supposed to ‘train’ itself on legitimately sourced content which an AI tool can convert into output, for our purposes text. LLMs range from simple, specialised tools, e.g. for use in marketing or customer services, to the very largest collections of datasets for general interrogation. In the latter case, this usually involves building extensive corpuses from sources where permission is generally not required, e.g. Wikipedia, government sites and the open web, but also pirated material.⁹

A study on ‘input’ sources has identified the following commonly used ‘pre-training’ corpora: web pages (open-source such as ‘CommonCrawl’, and ‘WebText’ Reddit posts); books and academic data (such as ‘Book Data’ general reading sources, ‘Project Gutenberg’ literary books, and ‘arXiv’ and ‘S2ORC’ academic papers and derivative datasets); Wikipedia; code (e.g. from ‘GitHub’ and ‘StackOverflow’); and a number of mixed data sources (such as ‘the Pile’ open source data set incorporating books, websites, codes, scientific papers and social media platforms).¹⁰ Another paper reported that OpenAI’s GPT-3 model¹¹ was trained with 499 billion ‘tokens’ of data derived from mostly online resources crawled on the web, including Wikipedia, the so-called ‘Book1’ dataset likely comprising Project Gutenberg’s public domain books, and a much larger ‘Book2’ dataset with unknown content.¹² This paper also reported widespread use of The Pile, but also included in its description the so-called ‘Book3’ dataset containing fiction and non-fiction books that appeared to be pirated.¹³ It further observed that, at the time of writing, OpenAI had not released much information on the training of GPT-4.

These datasets are consequences of the evolution of the information society, and were not constructed with AI in mind. Some datasets and web-based sources are intended to be ‘publicly accessible’ materials, even on an open access basis, but their exploitation in AI has had side-effects. For instance, the more indiscriminate the material trawled by an LLM, the greater the likelihood of so-called ‘hallucination’ in the production of output.¹⁴ Errors in content being interrogated, such as content derived from social media and networks and the open internet, can be replicated in outputs, for instance in texts generated and in key details such as citations. This has

⁹ Ethan Mollick. 2024. *Co-intelligence: Living and working with AI*. London: WH Allen. See page 33.

¹⁰ Wayne Xin Zhao et al. ‘A survey of large language models.’ *ArXiv*, 13 October 2024. See: <https://arxiv.org/pdf/2303.18223> (accessed 14 April 2025).

¹¹ GPT is the acronym for ‘Generative Pre-trained Transformer’, a type of large language model and framework for generative AI first introduced in 2018 by OpenAI. See: https://en.wikipedia.org/wiki/Generative_pre-trained_transformer (accessed 14 April 2025).

¹² Andres Guadamuz. ‘A scanner darkly: copyright liability and exceptions in artificial intelligence inputs and outputs’. *GRUR International*, Volume 73, Issue 2, February 2024, 111–127. See: <https://doi.org/10.1093/grurint/ikad140>.

¹³ Meta has allegedly used a pirated digital archive, a so-called ‘shadow library’ LibGen (Library Genesis), to train its AI models. See: Mark Sellman. ‘Politicians’ books pirated and ‘used to train Meta AI’’. *The Times*, 22 March 2025.

¹⁴ See Mollick at footnote 9, page 53.

particular ramifications for researchers. That said, there are reports of progressive improvements in both datasets and training processes, and that each iteration of major AI systems is reducing hallucinations and improving citation accuracy.

As AI develops, so does the availability of high quality datasets specifically aimed at deployment in LLMs. For example, the University of Cambridge's Polymathic AI initiative produces 'high quality' datasets to promote scientific discovery, some of which have already been released through the collaborative machine learning platform, Hugging Face, for others to use freely to develop their own AI models.¹⁵ In the legal profession, the Luminance LLM is trained on 150 million legal documents to assist in the analysis and the generation of draft contracts.¹⁶

So in sum, training of AI models is being conducted using content which is 'publicly accessible' in the widest sense of the term. In a complex data ecosystem, this content may be passively legitimate, paywalled,¹⁷ or actively open source. However, as we shall see later, 'scraping' for the purposes of training has been undertaken largely without content traceability being reported, let alone with permission. In the early stages of AI development, where the philosophy seems to have been to ask for forgiveness rather than permission, 'publicly accessible' material has been liberally exploited. This is now being challenged in courts by a wide range of right holders, raising fundamental questions on the ownership and exploitation of intellectual property rights.

1.3. AI tools

There are already a substantial number of generative AI tools in development or already on the market including Anthropic's Claude, DeepSeek's R1, Google's Gemini, Microsoft's Copilot, OpenAI's ChatGPT and xAI's Grok. These would be expected to be trained on at least the classes of resources outlined in the previous section, or equivalent. It will be interesting to see if any of these general models emerge or evolve as meaningful, sustainable tools for academic researchers. Or will more specifically developed and trained tools be needed, leveraging focussed, specialist resources of the type produced by Polymathic AI.

In the academic publishing world, access has evolved from aggregation and search of e-resources, through the development of reference and citation platforms, to early stage integration of generative AI functions. For instance, Scopus AI interrogates

¹⁵ University of Cambridge. 'New datasets will train AI models to think like scientists'. 2 December 2024. See:

<https://www.cam.ac.uk/research/news/new-datasets-will-train-ai-models-to-think-like-scientists> (accessed 14 April 2025).

¹⁶ John Gapper. 'Eye for small print helps AI reach deep into the legal industry'. *Financial Times*, 21 January 2025. See also: <https://www.luminance.com/> (accessed 14 April 2025).

¹⁷ See page 112 of Guadamuz at footnote 12.

metadata, abstracts and author profiles relating to Scopus contents, while ScienceDirect AI promises ‘meaningful and citable responses exclusively from millions of high-quality, peer-reviewed, full-text research articles and book chapters’.¹⁸ The Web of Science ‘Research Assistant’ product promises to help researchers find key papers faster using publication and citation data from its existing platform.¹⁹ Semantic Scholar advertises itself as a free, AI-driven search and discovery tool for academic papers sourced from a wide range of publishing and aggregator partners.²⁰

2. Ethical aspects

The rapid rise of AI has raised substantial ethical questions, which, with one exception, we will only briefly mention here. The concern with the greatest public profile is the potential for AI to surpass and outperform human intelligence and become uncontrollable and irreversible, the so-called ‘singularity’.²¹ A large group of signatories active in technology research have called for a pause on the development of systems more powerful than GPT-4 arguing that ‘Advanced AI could represent a profound change in the history of life on Earth’ with systems already ‘becoming human-competitive at human tasks’.²²

Another statement of AI practitioners drew attention to the extent of ‘data theft’ and use of cheap labour to develop for-profit tools, calling for regulation to enforce transparency on training and model architecture, and accountability for outputs.²³

Other criticisms highlight environmental impacts of AI through its consumption of enormous levels of natural resources, including energy and water coolants required to power the computing needs of data centres,²⁴ the expansion of transmission

¹⁸ Elsevier. ‘Scopus AI: Trusted content. Informed by responsible AI’. 2025 (<https://www.elsevier.com/products/scopus/scopus-ai>) and Elsevier. ‘Welcome to ScienceDirect AI. Eureka, every day’. 2025 (<https://www.elsevier.com/products/sciencedirect/sciencedirect-ai>) (both accessed 14 April 2025).

¹⁹ Clarivate. ‘Clarivate Launches Generative AI-Powered Web of Science Research Assistant’. 4 September 2024. See: <https://ir.clarivate.com/news-events/press-releases/news-details/2024/Clarivate-Launches-Generative-AI-Powered-Web-of-Science-Research-Assistant/default.aspx> (accessed 14 April 2025).

²⁰ See: <https://www.semanticscholar.org/about/publishers> (accessed 14 April 2025).

²¹ Tim Mucci. ‘What is the technological singularity?’. *IBM Think*, 7 June 2024. See: <https://www.ibm.com/think/topics/technological-singularity> (accessed 14 April 2025).

²² E.g. Yoshua Bengio et al. ‘Pause Giant AI Experiments: An Open Letter’, 22 March 2023. See: <https://futureoflife.org/open-letter/pause-giant-ai-experiments/> (accessed 14 April 2025).

²³ Timnit Gebru et al. ‘Statement from the listed authors of Stochastic Parrots on the “AI pause” letter’. DAIR Institute, 31 March 2023. See: <https://www.dair-institute.org/blog/letter-statement-March2023/> (accessed 14 April 2025).

²⁴ Mark Sellman, Adam Vaughan. “‘Thirsty’ ChatGPT uses four times more water than previously thought”. *The Times*, 4 October 2024.

networks and use of AI applications, as well as the challenges posed by ‘e-waste’ through turnover in computer servers and user devices.²⁵

The principal ethical consideration that will concern us here is the legal interest of right holders both with inputs into AI systems and outputs from them. The primary ethical concern with input is in uncontrolled use of copyrighted material to train AI models. At face value, existing copyright legislation might be expected to apply to the construction of LLMs. For instance, where content is used for input into training processes, even from open sources, there are formal matters to consider, such as observance of any restriction on commercial use, and a downstream obligation to cite in outputs. However, as we have seen above, there is substantial evidence that content is being widely used for commercial purposes without permission. Innovators are offering a variety of arguments that make a special case for the development of generative AI. One is roughly along the lines that the fungibility of content in LLMs does not threaten right holder interests individually, and that outputs cannot be meaningfully ascribed to individual works. A more technical justification is based on a wide interpretation of text and data mining legislation, a measure which was originally adopted for non-profit research purposes. A third approach proposes a re-ordering of the purpose of copyright in a new technological environment with changes to law to weight it in favour of continued AI innovation, for instance through the transitional use of regulated sandboxes permitting training of models in controlled environments. In effect, this is a plea that AI is a unique development in information technology for which existing rules should be adapted. One AI leader has even gone so far as to suggest that the power of AI technology is such that there will be ‘some change required to the social contract,’²⁶ and another making a call to ‘delete all IP law’.²⁷

These arguments are inventive, and some may indeed ultimately influence changes to legislation. However, as matters stand, the actions of the producers of many AI tools arguably contravene settled law. Individual and corporate owners of content are

²⁵ Alina Maria Vaduvam Kirk Chang. ‘A rising tide of e-waste, made worse by AI, threatens our health, the environment and the economy’. *The Conversation*, 29 November 2024. See: <https://theconversation.com/a-rising-tide-of-e-waste-made-worse-by-ai-threatens-our-health-the-environment-and-the-economy-244203> (accessed 14 April 2025).

²⁶ Falyn Stempler. ‘Panic after Sam Altman says AI will require ‘changes to social contract’ and society’. *Express US*, 26 January 2025. See: <https://www.the-express.com/tech/tech-news/161689/sam-altman-ai-societal-change-backlash> (accessed 14 April 2025). In a submission to a UK parliamentary science, innovation and technology committee, OpenAI called for a broad copyright exemption for AI, going so far as to reject a proposed ‘opt-out’ model intended to ameliorate the concerns of the creators of works. See: Georgia Lambert. ‘AI giants reject government’s approach to solving copyright row’. *The Times*, 3 April 2025.

²⁷ The quote is attributed to Jack Dorsey in a post of his on X dated 11 April 2025 (accessed 14 April 2025). See also: Mark Sellman. ‘Labour caving in on AI copyright laws, says Getty Images chairman’. *The Times*, 13 April 2025.

organising lobbying campaigns against laissez faire approaches,²⁸ as well as making direct challenges in courts, as will be discussed in more detail below.

AI outputs raise their own set of distinct issues. OpenAI acknowledges that as a user 'you (a) retain your ownership rights in Input and (b) own the Output'. They assign to the user 'all our right, title, and interest, if any, in and to Output'. Furthermore, ChatGPT requires users to ensure that their own inputs and the outputs generated from the tool ('content') do not violate any applicable law.²⁹ Clauses of this nature would appear at first sight to absolve AI companies of copyright concerns, also in light of contentions that the creation by LLMs of content for users is argued to be not human authorship of a work,³⁰ thereby falling outside of existing copyright law. However, users do have agency through their prompting of responses, and the issue of attribution of resources trained into LLMs has to be accounted for, and arguably even facilitated by AI producers.

3. Policy and regulation

Although there are already many models available globally for a variety of generative AI purposes, policy-makers and legislators in almost all affected jurisdictions are still grappling with AI's evolution. Official policy directions are in a state of flux, with general, and even conflicting, tendencies to encourage innovation, and to foster compromise between technological development and the property and exploitation interests of right holders. Other policy concerns address the protection of privacy and prevention of abusive outcomes through, for instance, the imposition of so-called 'guardrails'.

Attempts to set common principles and parameters at international level are also under way, although divisions have emerged between parties with regulatory and laissez faire approaches.³¹ To date, only the European Union (EU) has adopted comprehensive legislation through its so-called 'AI Act' of 2024.³² In the US, there is

²⁸ Mark Sellman. 'AI training on copyrighted works will 'wipe us out', say artists'. *The Times*, 17 December 2024; William Turvill. 'Elton John and Paul McCartney in harmony over the dangers of AI'. *The Times*, 26 January 2025.

²⁹ OpenAI. 'Europe Terms of Use'. [Effective: December 11, 2024]. See: <https://openai.com/en-GB/policies/terms-of-use/> (accessed 14 April 2025).

³⁰ Sercan Ozcan et al. 'ChatGPT: what the law says about who owns the copyright of AI-generated content', *The Conversation*, 17 April 2023. See: <https://theconversation.com/chatgpt-what-the-law-says-about-who-owns-the-copyright-of-ai-generated-content-200597> (accessed 14 April 2025).

³¹ Mark Sellman. 'UK and US snub France by refusing to sign AI summit declaration'. *The Times*, 11 February 2025.

³² Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence. *Official Journal of the European Union: L Series*, 12 July 2024. See: <http://data.europa.eu/eli/reg/2024/1689/oj> (accessed 14 April 2025).

no regulatory or policy guidance at federal level. AI companies are operating in a laissez faire environment with disputes being litigated into state level federal courts.

The UK is currently addressing AI legally through provisions on computer-generated content in its Copyright, Designs and Patents Act 1988,³³ with the prospect of more specific legislation such as an online safety act on illegal and harmful material due to come into force in 2025.³⁴ A government consultative process on AI concluded in late February 2025. This has given rise to divergences between competing interests, and may lead to further variations in policy and legislation.³⁵ There is also at least one significant legal challenge underway in a higher court based on claimed infringements of intellectual property rights.

Other potential legal avenues may emerge over time, particularly where entities exercise extensive market power, and where dominant players are perceived to abuse their positions, or where oligopolies cause distortions. In the past, the US and the EU have used antitrust powers to address anti-competitive practices, but actions of this nature tend to unfold over a long time frame. None are immediately in prospect, although in the US the White House issued an executive order in October 2023³⁶ where it addressed the promotion of competition.³⁷ More specifically, certain US state agencies³⁸ have been reported as paying attention to ‘algorithm-related collusion’.³⁹

3.1. The European Union - a legislative approach

The so-called ‘AI Act’ is an EU regulation comprising 180 recitals (introductory, contextual preamble), 113 legally applicable articles and 13 annexes, all creating a comprehensive framework to direct the development and oversight of AI in the Union. The EU’s policy approach is explicit in the recitals and states from the outset that AI should be a ‘human-centric technology’, serving as a ‘tool for people, with the ultimate aim of increasing human well-being’, with a ‘high level of protection of public

³³ See: <https://www.legislation.gov.uk/ukpga/1988/48/contents> (accessed 14 April 2025).

³⁴ Mark Sellman. ‘We won’t compromise online safety to appease Trump, minister signals’. *The Times*, 11 February 2025.

³⁵ Mark Sellman. ‘AI copyright shake-up could breach international law’. *The Times*, 7 March 2025.

³⁶ The White House. ‘Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence’, 30 October 2023. Available as ‘historical material’ at: <https://bidenwhitehouse.archives.gov/briefing-room/presidential-actions/2023/10/30/executive-order-on-the-safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence/> (accessed 14 April 2025).

³⁷ See section 5.3 of the order, including ‘addressing risks arising from concentrated control of key inputs, taking steps to stop unlawful collusion and prevent dominant firms from disadvantaging competitors, and working to provide new opportunities for small businesses and entrepreneurs’ and promoting ‘competition and innovation in the semiconductor industry’.

³⁸ The Department of Justice and the Federal Trade Commission.

³⁹ Alden Abbott, ‘Why Antitrust Regulators Are Focused On Problematic AI Algorithms’. *Forbes*, 13 March 2024.

interests as regards health, safety and fundamental rights' (paragraphs 6 and 7). It identifies classes of AI that require regulation, i.e. 'high risk AI' and 'general-purpose AI with systemic risk'.⁴⁰ It specifies concerns that must be protected, such as personal data, intellectual property rights, and identifies problems relating to biometrics and 'social scoring', and areas of activity where AI has particular impacts such as in education and employment. The regulation supports the promotion of innovation through nationally-deployed sandboxes as 'safe and controlled space for experimentation' (paragraph 138). The regulation intervenes on the matter of content used to train AI systems, emphasising the vital role of structured 'high-quality data' (paragraph 67) used transparently and legally (paragraph 107).

In the legally-binding text of the AI Act, Article 2 summarises the scope of the regulation. First of all it is intended to put in place EU-wide common rules for the placing on the market, putting into service and use of AI systems, including general-purpose models (see 'Chapter V' Articles 51 to 94). It explicitly prohibits certain AI practices impinging on personal privacy and integrity (see Article 5), and specifies requirements for high-risk AI systems and operators (see Articles 6 to 49). It introduces harmonised transparency rules for certain AI systems, and provides for official oversight. Finally, it governs measures to support innovation (Articles 57 to 63), primarily authorised 'sandboxes'.

Some of the stand-out provisions of the AI Act include the requirement for 'high-risk'⁴¹ AI systems to have models trained in accordance with quality criteria, including accounting for bias, special protection for personal data and specifications for input (training) data (see Articles 10 and 13).

General-purpose AI (such as large generative models) is treated in a similar manner to high-risk AI. Part of the control mechanism, which includes giving the European Commission a prominent, autonomous role, particularly through a specialist AI Office, features strict information provision obligations where 'systemic risk' is involved. For instance, generative AI models may be required to impart 'information on the data used for training, testing and validation, where applicable, including the type and provenance of data and curation methodologies' (see Annex XII, 2(c)).

⁴⁰ **General-purpose AI model:** 'an AI model, including where [...] trained with a large amount of data using self-supervision at scale, that displays significant generality and is capable of competently performing a wide range of distinct tasks regardless of the way the model is placed on the market and that can be integrated into a variety of downstream systems or applications [...]' (Article 3(63)).

Systemic risk: 'a risk that is specific to the high-impact capabilities of general-purpose AI models, having a significant impact on the Union market due to their reach, or due to actual or reasonably foreseeable negative effects on public health, safety, public security, fundamental rights, or the society as a whole, that can be propagated at scale across the value chain' (Article 3(65)).

⁴¹ See Article 6 definition and the classes listed in Annex III.

Article 53 requires providers of general-purpose AI models to keep up-to-date technical documentation which includes a detailed summary of content used for training.

Furthermore, providers are required to have a policy to comply with copyright and related rights, particularly with respect to text and data mining (TDM) relating to 'reproductions and extractions of lawfully accessible works'. One particular controversy relating to AI and TDM is the validity of so-called 'opt-out' mechanisms, also characterised as 'reservation of rights', which AI producers are invoking as a potential defence for data-scraping. If confirmed in legislation, these would place an onus on copyright holders to actively protect their interests through making reservation of rights statements, with failure to do so having consequences for the exploitation of their works.

Regulated sandboxes, which are intended to encourage innovation in real world environments for pre-marketing purposes, will have legal implications. When they enter into force from August 2026, they are formally expected to comply with personal data protection requirements. It is less clear how they will interact with copyright concerns.

Evolving law on training and 'opt-outs'

The Act lays out further elements of a regulated structure for the development of AI. Recital 105 of the Regulation considers the potential conflict between creators of content and its use in the training of generative AI models. It makes specific reference to TDM techniques as a common method for retrieving and analysing training content. This is an area where the EU has already established some ground rules. The Digital Single Market (DSM) directive of 2019⁴² provided for exceptions or limitations to copyright law permitting specific uses of TDM for automated computational analysis.⁴³ Article 3 of the directive accommodates the use of TDM for scientific purposes by research organisations (including universities and their libraries) and cultural heritage institutions for content to which they have lawful access. Article 4 enables a wider application of TDM 'to lawfully accessible works and other subject matter' except where right holders have expressly reserved the right to use the content concerned. This 'reservation of rights' is directly mentioned in Article 53.1(c) of the AI Act in the context of general-purpose models being obliged to draw up policies to comply with any TDM reserved rights.

⁴² Directive (EU) 2019/790 of the European Parliament and of the Council of 17 April 2019 on copyright and related rights in the Digital Single Market. *Official Journal of the European Union*. L 130/92, 17.5.2019. ELI: <http://data.europa.eu/eli/dir/2019/790/oj> (accessed 14 April 2025).

⁴³ Article 2(2) of the DSM directive describes TDM as 'any automated analytical technique aimed at analysing text and data in digital form in order to generate information which includes but is not limited to patterns, trends and correlations'.

This combination of the DSM directive and AI regulation suggests that there is a potential legal basis for training AI on content obtained by TDM. However, this is not yet clearly articulated in EU law, and there still appears to be a substantial difference between use of content for TDM and for AI purposes. This distinction has been considered in Germany in the Hamburg Regional Court in *Kneschke v LAION*.⁴⁴ This case concerned the use of TDM limited to the purpose of analysing image files for conformity with a pre-existing image description. In terms of AI, this was only about the preparatory⁴⁵ collection of training data and not the reproduction of the work during AI training. For the actual AI training, another download (i.e. reproduction) would be necessary. This legitimate use under TDM was not deemed to prejudice a stricter application of the three-step-test to the actual AI training.⁴⁶

Bearing in mind that the *Kneschke* case could be appealed as high as the European Court of Justice, and that there is some ambiguity in the relationship between the 2019 directive and the 2024 regulation (the 'Act'), for now producers of AI models are under notice via the Act to take at least the following steps to ameliorate their exposure under copyright law:

- Conduct training activities in line with a copyright law compliant policy under Article 53.1(c), whether explicitly for AI or TDM purposes, only on content for which no rights have been reserved by right holders; and
- Comply with Article 53 requirements on accounting for data used in training, and publicising detailed summaries of content used.

It should also be borne in mind that the formal exceptions or limitations to copyright law to date have been made for TDM purposes, and that if future case law decides there is a substantive difference between TDM and AI, full scale AI training may need a new legislative basis to use material that is not legally or freely available to producers.

⁴⁴ *Kneschke v LAION*, Judgment of Landgericht Hamburg, 27 September 2024 (Case no. 310 O 227/23). German version published at: <https://openjur.de/u/2495651.html> (accessed 14 April 2025).

⁴⁵ See Paul Goldstein et al. 'Kneschke vs. LAION – Landmark Ruling on TDM exceptions for AI training data – Part 1'. *Kluwer Copyright Blog*, 13 November 2024 at: <https://copyrightblog.kluweriplaw.com/2024/11/13/kneschke-vs-laion-landmark-ruling-on-tdm-exceptions-for-ai-training-data-part-1/> (accessed 14 April 2025).

⁴⁶ See the detailed analysis in Jonathan Pukas, Jan Bernd Nordemann. 'German Regional Court (Landgericht) of Hamburg paves the way to treat the reproduction of works as AI training data under the EU text and data mining exceptions'. *Kluwer Copyright Blog*, 25 October 2024 at: <https://copyrightblog.kluweriplaw.com/2024/10/25/german-regional-court-landgericht-of-hamburg-paves-the-way-to-treat-the-reproduction-of-works-as-ai-training-data-under-the-eu-text-and-data-mining-exceptions/> (accessed 14 April 2025).

3.2. The United States - litigating change?

At US federal level, attempts by the Biden Administration to introduce federal legislation in early 2023⁴⁷ bore no fruit. However, the White House issued a detailed executive order in late 2023⁴⁸ with a similar level of detail to the EU Act preamble, but which was not legally enforceable. The order set out eight guiding principles and priorities in its section 2, including being 'safe and secure', 'responsible', in line with 'privacy and civil liberties' and leading to 'technological progress'. However, these guides have been jettisoned by the Trump Administration, and a more laissez faire approach was signalled with the appointment of an AI chief who favours technology industry calls for a policy 'friendlier' to innovation.⁴⁹

At state level, an AI bill in California which had sought to impose safety tests and enable shut down by humans was successfully quashed.⁵⁰

Despite the lack of legislative guidance, conflicts are emerging in a wide range of legal challenges, mostly in federal courts in California and New York. They have been instigated by a variety of content type creators or owners including for books,⁵¹ news,⁵² images⁵³ and open-source software code,⁵⁴ primarily concerning the use of their copyrighted material for LLM training.

⁴⁷ Ryan Tracy. 'Biden administration weighs possible rules for AI tools like ChatGPT'. *Wall Street Journal*, 11 April 2023.

⁴⁸ The White House. 'Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence'. 30 October 2023. Available as 'historical material' at: <https://bidenwhitehouse.archives.gov/briefing-room/presidential-actions/2023/10/30/executive-order-on-the-safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence/> (accessed 14 April 2025).

⁴⁹ Preetika Rana. 'Trump plans to appoint Musk confidant David Sacks as AI, crypto czar'. *Wall Street Journal*, 5 December 2024.

⁵⁰ See Rana at footnote 49; and Preetika Rana. 'AI companies fight to stop California safety rules'. *Wall Street Journal*, 7 August 2024.

⁵¹ *Tremblay v OpenAI Inc.*, District Court, N.D. California, Case 3:23-cv-03223; *Author's Guild v OpenAI*, District Court, S.D. New York, Case 1:23-cv-08292; *Kadrey v Meta Platforms*, District Court, N.D. California, Case 3:23-cv-03417.

⁵² *Wall Street Journal and New York Post v Perplexity AI Inc.*, District Court, S.D. New York, Case 1:24-cv-07984; *New York Times v. Microsoft Corp.*, District Court, S.D. New York, Case 1:23-cv-11195; *Raw Story Media v OpenAI Inc.*, District Court, S.D. New York, Case 1:24-cv-01514-CM

⁵³ *PM v OpenAI LP*, District Court, N.D. California, Case 3:2023cv03199; *Anderson v Stability AI*, U.S. District Court, N.D. California, Case 3:23-cv-00201; *Getty Images Inc. v Stability AI*, District Court, D. Delaware, Case 1:23-cv-00135.

⁵⁴ *DOE v GitHub Inc.*, District Court, N.D. California, Case 4:22-cv-06823. See Pamela Samuelson. 'Generative AI meets copyright'. *Science*, 14 July 2023 (Vol. 381, issue 6654, p158).

Case law decisions

The one case determined to date is *Raw Story Media v OpenAI Inc.* decided in November 2024.⁵⁵ Raw Story claimed thousands of their copyrighted works of journalism were ‘scraped’ and used in OpenAI’s training sets, that copyright management information (CMI), i.e. identifying and rights metadata, was removed beforehand, and that this was a violation meriting damages. The legal provision concerned requires ‘concrete injury’. The court denied such injury on the basis that Raw Story had not identified a query made on ChatGPT that disseminated their work or alleged any actual adverse effects.

An additional request was made for an injunction to remove all works from which CMI was extracted because of the risk that ChatGPT would reproduce work without the attendant CMI. The court rejected the injunction arguing that the chance that ChatGPT would plagiarise actual Raw Story content seemed remote, rationalising that ChatGPT included massive amounts of information from innumerable sources.

The court left it open to Raw Story to return with an amended complaint, but expressed doubts as to whether they could identify a relevant injury, although it held out slightly better prospects for further pleading on an injunction.

Themes in undecided cases

Book authors have taken two class actions against OpenAI on the use without consent of their copyrighted works to train ChatGPT. The *Tremblay* action filed in June 2023 alleged such a use and the generating of profit therefrom. When training GPT-3, OpenAI disclosed that 15% of the training dataset came from “two internet-based books corpora” (“Books1” and “Books2”) without revealing which titles of which they are comprised. It was submitted that the dataset contained approximately 350,000 titles, again unidentified. The plaintiffs contend their works form part of this dataset. The Authors Guild alleged mass theft of their fiction works which represent their creative literary expression, and the output of derivative works through outputs from LLM, all without payment of a reasonable licensing fee.⁵⁶

⁵⁵ See footnote 52. Much of the detail for these cases has been obtained from documents (primarily complaints) available as downloadable PDFs through the free online service Court Listener (<https://www.courtlistener.com/>). Search for cases on the home page, or enter the unique case number in a search engine and navigate accordingly.

⁵⁶ See *Tremblay* and *Authors Guild* cases at footnote 51. The Authors Guild itself has reported on a licensing deal proposed by a publisher to permit use of non-fiction books for AI purposes subject to express permission and a fee per title. The Guild criticised the proposed 50/50 fee split between publisher and author as being against the latter’s rights, and also noted limits on the extent of a book’s text that can be used in outputs. See: ‘HarperCollins AI Licensing Deal’. *The Authors Guild*, 19 November 2024 at <https://authorsguild.org/news/harpercollins-ai-licensing-deal/> (accessed 14 April 2025); and Ella Creamer. ‘HarperCollins to allow tech firms to use its books to train AI models’. *The Guardian*, 19 November 2025.

At least two cases have been taken by newspapers. The first was instigated by the New York Times (NYT) against OpenAI in May 2024.⁵⁷ The newspaper alleged that GPT-4 could be used to generate several paragraphs of Times content as outputs in response to user prompts. OpenAI claimed the NYT went to great effort with multiple prompting to achieve these results. The NYT argued this type of use of AI imperils journalism, while OpenAI highlighted the benefits of ‘normal use’ of generative AI.

The second was lodged in October 2024.⁵⁸ The Wall Street Journal and New York Post have taken a case against Perplexity AI alleging a massive amount of illegal copying of publishers’ copyrighted works and diverting customers and critical revenues away from the copyright holders. The publishers have emphasised the length to which they go to produce content and its consequent value. They have alleged distinct infringements relating to inputs into Perplexity’s AI product and its output ‘answers’, and the impact on their reputations, looking for reliefs on all these counts.

A class action taken against Stability AI in January 2023 alleged the use of billions of copyrighted images without permission or compensation, both for training and output purposes. The plaintiffs argued that derived images compete in the marketplace with the original images, not least those generated ‘in the style’ of a given artist. Getty Images also took a case against Stability AI in February 2023 for the use without permission of 12 million of its images to ‘build a competing business’.⁵⁹

The plaintiffs in *PM v OpenAI LP* alleged in June 2023 that generative AI models collect, store, track, share, and disclose private information of millions of users of, for instance, social media, and music and video applications. They claimed OpenAI uses stolen private information, including personally identifiable information, from hundreds of millions of internet users, including children of all ages, without their informed consent or knowledge

Doe v GitHub Inc. is another class action, this time taken by owners of copyright interests in materials made available publicly on GitHub.⁶⁰ Although originally open source, GitHub was acquired by Microsoft in 2016 and ultimately ended up in 2022 as part of the Copilot product designed to assist software coders by providing or filling in blocks of code using AI. The complaint claimed that the ‘Defendants have been cagey about what data was used to train the AI’, but ‘they have conceded that the training data includes data in vast numbers of publicly accessible repositories on GitHub, which include and are limited by Licenses’. Furthermore, the ‘Defendants stripped Plaintiffs’ and the Class’s attribution, copyright notice, and license terms

⁵⁷ See footnote 52.

⁵⁸ Alexandra Bruell, ‘Wall Street Journal, New York Post sue AI startup Perplexity, alleging ‘massive freeriding’, *Wall Street Journal*, 22 October 2024. See also footnote 52.

⁵⁹ See *Anderson v Stability AI* and *Getty Images v Stability AI* in footnote 53.

⁶⁰ *Getty Images v Stability AI, Ltd.* High Court of England and Wales. Case No: IL-2023-000007.

from their code in violation of [...] Licenses and [...] rights'. The defendants then used Copilot to distribute anonymized code to users as if created by Copilot. The plaintiffs proceeded to claim for damages and injunctive relief.

The issue of the removal of copyright management information from works prior to training, which was unsuccessful in *Raw Story Media v OpenAI Inc.* above, got more traction in a case taken by book authors against Meta for the production of the Llama models. The judge in *Kadrey et al vs Meta Platforms* has allowed the plaintiffs to proceed with the argument that Meta's removal of copyright notices violates the Digital Millennium Copyright Act.⁶¹

It should be noted that the above scenarios are mainly allegations by plaintiffs with some counter-arguments by defendants. Final court determinations will likely favour a narrow range of facts and claims. However, the following broad trend can be identified among copyright holders across a broad range of content types where all are contesting the use of their works for AI training without permission and for the production of output that replicates or resembles their work on some level without attribution. More broadly, users of social media and entertainment applications are objecting to the use of information relating to their use of platforms on privacy grounds. In almost all cases, plaintiffs are looking for compensation and some form of injunctive relief. This suggests some plaintiffs see their issue as an economic one, while others argue that their content needs to be treated in line with essential rights.

Some academic analysis

The undecided cases discussed above may ultimately be resolved on traditional copyright first principles. However, there have also been academic enquiries into how the nature of AI might challenge how these principles are understood. A complex consideration of copyright aspects of AI is laid out in an article by Fenwick and Jurcys,⁶² in this case in point concerning the use of AI for music production purposes. It involves the use of generative AI tools by the dance music producer, David Guetta, to 'Write a verse in the style of Eminem about the Future Rave'.

While re-mixing and sampling has long been a feature of the music industry, and formally regularised through permission and royalty arrangements, the use of AI to similar ends has introduced entirely new factors, such as interrogating indeterminate content 'trained' from 'publicly available' sources in large language models to produce output which does not attribute or credit trained material.

⁶¹ See footnote 51. See also: Thomas Claburn. 'Judge says Meta must defend claim it stripped copyright info from Llama's training fodder'. *The Register*, 11 March 2025, available at: https://www.theregister.com/2025/03/11/meta_dmca_copyright_removal_case/?td=rt-3a (accessed 14 April 2025); and Mark Sellman. 'Politicians' books pirated and 'used to train Meta AI'. *The Times*, 22 March 2025.

⁶² Mark Fenwick, Paulius Jurcys. 'Originality and the future of copyright in an age of generative IT', *Computer Law & Security Review: The International Journal of Technology Law and Practice*. (Vol. 51, 2023). See: <https://doi.org/10.1016/j.clsr.2023.105892>.

Fenwick and Jurcys examined fundamental aspects of US copyright law to assess the potential legality of the process Guetta used to produce his music clip. Among the general copyright and regulatory issues raised are:

- The 'legality of data scraping and using publicly available information' both copyrighted and 'in the public domain' to train machine learning and AI models. Are permissions or licenses required to conduct such scraping, and is the use of scraped data a copyright infringement in itself;
- Data privacy and image rights;
- Transparency and oversight for the data used to train LLMs, versus, for instance, AI content as a 'black box'; and
- The legal status of user outputs from AI tools, i.e. as to whether they are 'original' or 'derivative'.

In essence, the authors argued that the clip was the outcome of a 'creative use' of generative AI in a contemporary context. This was on the basis of Guetta prompting an AI tool and its generation of output for one element of his composition, his use of a similar process in another tool for another element, re-prompting and refining, followed by integrating of the two outputs to produce his final output. They saw this multi-stage process of a network of human, corporate and machine 'actors' as a new vision of creativity.

They proceeded to consider whether the clip passed the 'originality' test in US copyright law, noting that 'expressions' not 'ideas' are protected. Ideas are intangible, should be freely used to promote scientific development and progress, and constitute free speech. It is the formal expression which is copyrightable. Originality is determined by being created independently. The level of creativity required is 'modest', a low threshold. The independence criterion requires work to emanate from a natural person even if a tool is used, i.e. there must be a human in the loop. The authors ultimately concluded that the production of the clip was probably original and creative enough to attract the protection of copyright law.

3.3. The United Kingdom - legal evolution

Policy, regulation and conflicting interests

The UK is in a period of transition on AI regulation. There is a holding view that the Copyright, Designs and Patents Act 1988, which protects 'computer-generated' works, covers AI-generated work.⁶³ However it is not clear how this governs, e.g., use of data for training without licence, or how it addresses the conflict between

⁶³ Matthew Sparkes. 'AI copyright'. *New Scientist*, 8 October 2022 (Vol. 256, issue 3407, p17). See: [https://doi.org/10.1016/S0262-4079\(22\)01807-3](https://doi.org/10.1016/S0262-4079(22)01807-3).

creators of works and commercial generative AI practices. Accordingly, successive governments have been attempting to find solutions to plug the gap.

The Conservative government undertook a review of the regulation of digital technologies, including AI, which resulted in the so-called Vallance Report of 2023.⁶⁴ This report observed that the UK had a short window to become a top location for ‘foundational AI companies’, and that other countries had been quicker to provide a ‘friendly regulatory environment for innovators’. Recognising the various levels of acceptable risk, the report recommended a multi-regulator AI sandbox for a six-month period where rules would be relaxed and experimentation encouraged under supervision. The government was invited to announce a clear policy on the relationship between intellectual property law and generative AI, with clarification around the role of text and data mining mechanisms deployed by AI firms. Barriers to the use of ‘publicly available’ copyright and database materials were also questioned. In addition, concerns were raised over privacy aspects concerning ‘public data’.

The UK government accepted the recommendations of the Vallance Report,⁶⁵ including agreeing to (i) the establishment of the sandbox, (ii) the production of a code of practice by the Intellectual Property Office (IPO) on access to copyrighted work for input and protections on outputs in the context of ‘reasonable’ licenses from right holders and (iii) facilitating greater industry access to public data. By late 2024, it had been reported that attempts by the IPO to find legal middle ground between technology and creative industries had reached an impasse. However, the new Prime Minister is reported as having said that news organisations will be paid for allowing AI to use their work,⁶⁶ with a solution to related copyright disputes to be proposed by the end of 2024⁶⁷ and legislation adopted in the second half of 2025.⁶⁸ Part of the difficulty to be resolved concerns the conflicting views of an AI industry which is lobbying for ‘liberalisation’ of copyright law and a creative sector, including major media and entertainment corporations, which objects to any attempt to ‘degrade’ current law. Mooted solutions include mandatory disclosure of ‘scraped’ source material, ‘opt-out’ clauses (‘rights reservation’), and an express permission

⁶⁴ HM Government. ‘Pro-innovation regulation of technologies review: Digital technologies’. March 2023. See:

https://assets.publishing.service.gov.uk/media/64118f0f8fa8f555779ab001/Pro-innovation_Regulation_of_Technologies_Review_-_Digital_Technologies_report.pdf (accessed 14 April 2025).

⁶⁵ HM Government. ‘HM Government response to Sir Patrick Vallance’s pro-innovation regulation of technologies review: Digital technologies’. March 2023. See:

https://assets.publishing.service.gov.uk/media/6410aa2ce90e076cc6e370ef/HMG_response_to_SPV_Digital_Tech_final.pdf (accessed 14 April 2025).

⁶⁶ Chris Smyth, Mark Sellman. ‘Newspapers must be paid if AI uses their archive, Sir Keir Starmer says’. *The Times*. 28 October 2024.

⁶⁷ Tom Saunders. ‘New law may be needed to end AI copyright disputes.’ *The Times*. 2 October 2024.

⁶⁸ Oliver Wright, Mark Sellman. ‘Celebrities get new right to fight against AI piracy.’ *The Times*. 14 December 2024.

requirement.⁶⁹ The UK government opened a consultation process in late 2024 which concluded in late February 2025.⁷⁰ The official position in this consultation leans towards a copyright exception in favour of commercial AI training, leaving it to copyright holders to rely on a rights reservation mechanism with the possibility of compensation through licensing. This has met with stiff resistance from legal experts,⁷¹ publishers⁷² and content creators themselves.⁷³ The government is expected to issue its response to the consultation in the summer or September.⁷⁴

Case law

Litigation similar to that under way in the US is before UK courts. Getty Images has taken a case against Stability AI. The complaint is that the Stability AI has 'scraped' millions of images from Getty Images websites without consent, and used those images unlawfully to train and develop their Stable Diffusion product, and also that the product's output of synthetic images is an infringement in that it reproduces a substantial part of copyright works.⁷⁵ The case is expected to go to trial in the High Court in summer 2025.⁷⁶

In another instance, it has been reported by a UK source that stock photograph agencies Getty Images and Shutterstock have taken unilateral practical action removing images from their platforms to prevent them being used by AI models like DALL-E, Stable Diffusion and Midjourney.⁷⁷

3.4. General legal directions of travel

At this point it is difficult to draw clear conclusions from the various scenarios unfolding in the EU, US and UK. In commercial terms, of the three jurisdictions, the

⁶⁹ See Wright, Sellman at footnote 68.

⁷⁰ Intellectual Property Office et al. 'UK consults on proposals to give creative industries and AI developers clarity over copyright laws'. 17 December 2024. See: <https://www.gov.uk/government/news/uk-consults-on-proposals-to-give-creative-industries-and-ai-developers-clarity-over-copyright-laws> (accessed 14 April 2025).

⁷¹ Mark Sellman. 'Copyright law reforms may face review over breaches'. *The Times*, 28 February 2025.

⁷² Mark Sellman, 'AI copyright shake-up could breach international law'. *The Times*, 7 March 2025.

⁷³ Mark Sellman. 'AI copyright law plan 'idiotic', says Sir Cameron Mackintosh'. *The Times*, 5 March 2025.

⁷⁴ Mark Sellman. 'Rollback on AI copyright could lead to parliamentary "ping-pong"'. *The Times*, 1 April 2025.

⁷⁵ Summarised from paragraph 8 of a ruling of the High Court of 1 December 2023 in a failed application to have Getty Images'. See: https://archive.org/stream/getty-v-stability-ai-11202023-sjruling-uk/GettyvStabilityAI11202023SJRulingUK_djvu.txt (accessed 14 April 2025).

⁷⁶ Cerys Wyn Davis. 'Getty Images v Stability AI: the implications for UK copyright law and licensing'. *Pinsent Masons: Out-Law / Your Daily Need-To-Know*. 29 April 2024 (see: <https://www.pinsentmasons.com/out-law/analysis/getty-images-v-stability-ai-implications-copyright-law-licensing> (accessed 14 April 2025)).

⁷⁷ See Sparkes at footnote 63.

AI industry is most active in the US. In regulatory terms, the EU has the most comprehensive framework in place. Leaving aside the merits of business interests, property rights, and the public interest, a number of developments should be closely monitored in the course of 2025, including:

- The resolution of litigation in the US. These are likely to unfold from mid-2025, with focus on whether trends favour vindication of property rights, compromise based on compensation, or unconstrained use by AI enterprise;
- The interpretation of 'publicly accessible'. This is a term understood permissively by AI advocates to include even material behind paywalls. Apart from the primary impact on exploitation rights, the practice of stripping away CMI/metadata also militates against accurate citation in chatbot outputs. Courts may rule on this particular aspect;
- The UK's approach to conflicting stances of content producers and AI enterprise. The former are arguing that the very economic basis of their output is being undermined by current AI practices. They also argue that fundamental property rights established internationally under the Berne Convention are being contravened both by AI industry practices and government proposals to impose 'opt-out' obligations on right holders;
- The EU position on the 'opt-out'/reservation of rights' mechanisms. The AI Act is not due to come fully into force until 2026, and, although it is a comprehensive piece of legislation, many practical elements need to be brought into play. It remains to be seen whether its regulatory approach discourages innovation within the EU, or makes doing AI business attractive and sustainable from a compliance perspective.

Up to now, AI producers have been looking for forgiveness rather than permission, all the while investing heavily in building their LLMs and chatbots. This has led to right holders both to take direct action in courts and to lobby legislators to protect the value of their property. If the official tendency is towards vindicating property rights, this will likely add to the future cost of developing LLMs. If, on the other hand, technological innovation is favoured, this will constitute a fundamental change to the exploitation of copyrighted works.

4. Implications for higher education and research libraries

4.1. The higher education sector

Generative AI is now widely used in education and research and is influencing teaching, learning and assessment. Many diverging views have been advanced on its impacts and efficacy, and the jury is still out on definitive outcomes. On the one hand, tools like ChatGPT are argued to have the potential to enhance student learning, and on the other they are perceived to unduly distort assignment writing in a manner some describe as bad practice and even misconduct.

Some emerging trends

Some institutions, at least in the early stages of the marketing of ChatGPT and similar tools, adopted an approach of banning their 'unauthorised use'. For instance, a survey of student use of ChatGPT in the United States conducted in August and September 2023 reported that 29% of their institutions prohibited ChatGPT, with only 7% permitting it.⁷⁸

A plethora of challenging impacts, such as deficiencies in output, and plagiarism and its management, has led educators to rethink teaching and assessment design, in particular online assessment.⁷⁹ Other commentators have gone so far as to project the use of AI as a threat to higher education, as it engenders a tendency toward simple provision of answers to queries as opposed to fostering cognitive development.⁸⁰ Some institutions such as Brown, Dartmouth, the Institut d'études politiques de Paris (SciencePo) and Yale are reintroducing standardised written exam testing and even entrance exams to counter the use of ChatGPT to (over)prepare written assignments.⁸¹

In contrast, an international panel of HE educators and experts produced a more comprehensive and methodical assessment of AI in a study published in 2024.⁸² As an exercise intended to raise awareness and counsel a cautious, iterative approach,

⁷⁸ Clare Baek et al. "ChatGPT seems too good to be true": College students' use and perceptions of generative AI'. *Computers and Education: Artificial Intelligence*, Volume 7, December 2024. See: <https://doi.org/10.1016/j.caeai.2024.100294>.

⁷⁹ Dirk H.R. Spennemann et al. 'ChatGPT giving advice on how to cheat in university assignments: how workable are its suggestions?'. *Interactive Technology and Smart Education*, Volume 21, Issue 4, Pages 690 - 707, 30 October 2024.

⁸⁰ Anthony Seldon. 'Warning from AI is stark: we have two years to save learning'. *The Times*, 25 October 2024.

⁸¹ Maurício Alcenaar. 'Elite French university revives entrance exam to combat AI fears'. *The Times*, 23 October 2024.

⁸² Aras Bozkurt et al. 'The manifesto for teaching and learning in a time of generative AI: A critical collective stance to better navigate the future.' *Open Praxis*, 16(4), pp. 487–513. (2024). See: <https://doi.org/10.55982/openpraxis.16.4.777>.

the paper identified a range of positive⁸³ and negative impacts⁸⁴ of generative AI on teaching and learning. Without claiming to issue prescriptive findings, it identified 35 individual themes across the positive and negative categories, in each case summarising some key headline considerations, complemented by a text passage with more developed insight. It is clear from this study that its authors perceive a widespread use of generative AI, and are building the case for its potential, properly deployed, to enhance 'human agency' and 'ethical responsibility' in higher education.

Impact on students

In practical terms, as far back as early 2023, a survey found that 30% of US college students were already using ChatGPT to complete written homework assignments, and that close to 60 percent had used it on more than half of their work.⁸⁵ Data from later in 2023 reported that 33% of college students used ChatGPT monthly for general and writing tasks, although 67% never used it for computer programming.⁸⁶ A 2024 UK study of Russell Group students revealed almost universal recourse to chatbots, with almost 20% copying directly from outputs, half of these submitting chatbot output as complete pieces of coursework, and, incidentally, as few as 1 in 400 being punished for perceived AI misuse.⁸⁷

At a Stanford University event, various views were offered on AI's impact on education and research.⁸⁸ For instance, would writing's role in the development of critical thinking be undermined? Others suggested that it was no longer necessary to write all essays without AI support, and that this in fact might actually encourage a greater emphasis on editing and curating, forcing a deeper engagement overall. In early 2023, a researcher of emerging technologies at Harvard counselled on AI's limitations, characterising ChatGPT as 'the next step beyond a search engine' giving you what it thinks you want, as opposed to simple lists of results, while acknowledging that a transition was underway from formal essay composition, which AI can perform 'beautifully', to producing narratives with a human, creative dimension, which AI cannot.⁸⁹

⁸³ Such as efficiency, personalised learning, preparing for work, and innovation.

⁸⁴ Such as the digital divide, bias, ethics, and integrity.

⁸⁵ 'Study: 30% of College Students Have Used ChatGPT for Essays'. *Government Technology*, 25 January 2023. See:

<https://www.govtech.com/education/higher-ed/study-30-of-college-students-have-used-chatgpt-for-essays> (accessed 14 April 2025).

⁸⁶ See Baek et al. at footnote 78.

⁸⁷ Fintan Hogan. 'AI cheats 'slip under radar' as few university students penalised'. *The Times*, 6 April 2025.

⁸⁸ Claire Chen. 'AI Will Transform Teaching and Learning. Let's Get it Right'. *HAI: Human-Centered Artificial Intelligence*, 9 March 2023. See:

<https://hai.stanford.edu/news/ai-will-transform-teaching-and-learning-lets-get-it-right> (accessed 14 April 2025).

⁸⁹ Jill Anderson. 'Educating in a World of Artificial Intelligence: Chris Dede discusses how education can evolve to work with - rather than fight against - artificial intelligence'. Transcript of *Harvard Graduate School of Education Edcast*, 9 February 2023. See:

In contrast, students surveyed on their use of ChatGPT expressed concern about losing their creativity, particularly among non-STEM fields where personal voice and creativity are more central. STEM majors were more concerned about lack of confidence 'in the accuracy of information provided by ChatGPT'.⁹⁰

Institutional guidance

Despite concerns, the spread of AI continues unabated. Ultimately, it is a matter for institutions to decide whether they wish to endorse the use of AI or not. The Russell Group of Universities in the UK has committed to the 'ethical and responsible use of generative AI'.⁹¹ The Group wants to ensure that generative AI tools are used to enhance teaching practices and student learning experiences, particularly so students develop future skills and educators develop innovative methods of teaching. Member universities are to support student AI literacy, equip staff to use AI tools effectively and appropriately, while ensuring academic rigour and integrity. While risks are identified, such as privacy, bias and inaccuracy, AI's positive potential is recognised. As AI use varies within disciplines, academic departments are being encouraged to produce local subject-relevant guidance embedded in institutional policy, and dialogue between staff and students is also encouraged. The principles acknowledge that AI technologies are developing and that new generative tools will become available, further noting that some of these tools may be subscriber-based or behind paywalls.

In the EU, the climate on using generative AI is significantly determined by the AI Act. In principle, human-centric and trustworthy AI must ensure health, safety and fundamental rights protecting against harmful effects. Generative AI tools produced in compliance with these terms are likely to have bridged key ethical concerns, and this will go a long way to reassuring institutions on their deployment by staff, students and researchers. Furthermore, where AI is used in a work context, Article 4 of the Act requires employers to ensure that relevant staff are AI literate. In the educational context, it would seem that the primary considerations concern academic integrity, the efficaciousness of any tools used, and whether or how they might be best integrated into teaching and learning, assessment and research practice.

In Ireland, the National Academic Integrity Network⁹² has published guidelines for educators on AI which recognise the evolving nature of generative AI, but which recommend practical approaches for the various stakeholders in higher education

<https://www.gse.harvard.edu/ideas/edcast/23/02/educating-world-artificial-intelligence> (accessed 14 April 2025).

⁹⁰ See Baek et al. at footnote 78.

⁹¹ Russell Group principles on the use of generative AI tools in education. See: https://russellgroup.ac.uk/media/6137/rg_ai_principles-final.pdf (accessed 14 April 2025).

⁹² Established in 2019 by Quality and Qualifications Ireland.

including useful context for users.⁹³ The section on what students need to know is an appraisal of expectations to be placed both on users and on the tools of which they may be trying to make sense.

CONUL institutions themselves, and specialised units within them, are increasingly laying down guidelines on the use of AI.⁹⁴ In general, they are informed by mandatory requirements in the EU's AI Act, such as the requirement to provide staff with appropriate training. These guidelines will no doubt evolve as regulatory impacts emerge and as institution-wide effects on teaching, learning and research come into clearer focus.

4.2. AI products in education and research

Consumer-oriented general-purpose chatbots such as ChatGPT are based on huge volumes of 'publicly available' content and can produce versatile outputs, albeit with the risk of considerable inaccuracy in generated texts and in citations. However, these are only one type of tool. Other elements relevant to generative AI range from high quality curated datasets shared to improve LLM training (such as Polymathic AI), through content-rich publisher sources adapted by AI functions (such as Scopus AI), potentially to chatbot-type products designed to generate substantial output for educational consumers.

Already, some of the research literature heretofore searchable through databases and discovery platforms is being augmented by AI functions. The principal output at this point seems to be limited primarily to the production of summaries and citations from abstracts,⁹⁵ derived from limited runs of holdings.⁹⁶ However, the next level of functionality will involve using a more complete range of available published sources.

⁹³ See National Academic Integrity Network. 'Generative artificial intelligence: Guidelines for educators'. 1st edition, August 2023. Available at: <https://www.qqi.ie/sites/default/files/2023-09/NAIN%20Generative%20AI%20Guidelines%20for%20Educators%202023.pdf> (accessed 14 April 2025).

⁹⁴ See various policies, statements and FAQs from the following institutions: [DCU](#); [MU](#); [QUB](#); [TCD](#); [UCC](#); [UCD](#); [UL](#) and [UU](#) (accessed 14 April 2025). Please note that these links are indicative and do not purport to record official policies. Furthermore, as institutional responses develop over time, original texts and URLs are likely to vary.

⁹⁵ See Aster Zhao. 'Trust in AI: Evaluating Scite, Elicit, Consensus, and Scopus AI for Generating Literature Reviews', 20 March 2024. See: <https://library.hkust.edu.hk/sc/trust-ai-lit-rev/> (accessed 14 April 2025).

⁹⁶ Scopus AI currently limits its summaries function to documents from 2003 on. See: 'Which Scopus content does Scopus AI draw on?' FAQ at: <https://www.elsevier.com/products/scopus/scopus-ai> (accessed 14 April 2025). ScienceDirect AI interrogates large volumes of full-texts, but without specifying date ranges. See: Elsevier. 'Welcome to ScienceDirect AI. Eureka, every day'. 2025 (<https://www.elsevier.com/products/sciencedirect/sciencedirect-ai>) (accessed 14 April 2025).

In this regard, for instance, one observer has noted that the real value in a platform like Scopus AI lies in being able to access the full text of articles.⁹⁷

Researchers themselves may query whether the greater benefit to them is the use of AI to generate reliable relationships between descriptions of works (metadata) in order to augment the retrieval of relevant sources for their own appraisal, to produce summaries of items of interest, or to derive 'creative' outputs from the full texts themselves. Herein lies a key question on the purpose of AI for research - does it simplify and deepen the iterative research process, or does it (or can it) produce text of value to the process of composition?

From a research output perspective, is there a sustainable level of additional productivity to be gained over time in the mining of existing information? And will the generative delivery of content be reflected in substantial advancement of knowledge? Or does scholarly advancement still effectively depend on close engagement with subject corpuses? As Spennemann et al. noted, generative AI has been used, frequently and with considerable success, as a brainstorming tool to collate and summarise information.⁹⁸ However, unless there are real advances beyond the current iteration of, say GPT-4, or unless the LLM approach on which GPT is based is replaced, perhaps the most productive use of AI is to make greater sense of existing specialist datasets.

As has been noted with the larger generative AI tools, training and testing is an expensive process, and there are still problems with hallucinations and citation errors. Accordingly, it remains to be seen whether academic platforms have the resources to be able to construct AI that can effectively interrogate source materials and compose outputs.

4.3. AI and research libraries

From a library practice perspective, literacy skills programmes and tools should consider what generative AI can realistically achieve. Answers may prove to be elusive, but this should at least include an honest appraisal of available search tools, and guidance on ethics, including reinforcing conventional concerns such as citing and referencing. Much practical direction has already been provided by CONUL libraries in specific AI guides, or integrated into existing guides.⁹⁹

⁹⁷ Teresa Kubacka. 'Guest Post - There is More to Reliable Chatbots than Providing Scientific References: The Case of ScopusAI'. *The Scholarly Kitchen*, 21 February 2024. See: <https://scholarlykitchen.sspnet.org/2024/02/21/guest-post-there-is-more-to-reliable-chatbots-than-providing-scientific-references-the-case-of-scopusai/> (accessed 14 April 2025).

⁹⁸ Spenneman (see footnote 79), page 704.

⁹⁹ See libguides focused primarily on AI at: [QUB](#); [RCSI](#); [TCD](#); [TUD](#); [UCD](#) and [UG](#) (accessed 14 April 2025). Please note that these links are indicative and do not purport to officially or permanently record

Widely used scholarly publishing platforms have AI components, but, at least for now, act more like search engines, albeit with augmented features. However, as publishers and aggregators improve their generative AI capabilities, these functions may improve, and even emerge in time as new subscriber-based platforms, which librarians will need to evaluate for suitability for their user communities. It is likely that in academic institutions, library literacy engagement will be complementary to guidance provided by lecturers. However, librarians in particular can continue to leverage their expertise in helping identify and evaluate specialist and high quality information sources, and their relevance to AI tools.

Implications for AI-generated outputs in educational and research contexts are worth a particular mention. OpenAI vests copyright for ChatGPT outputs in the user. However, users should be aware that this may not necessarily absolve them of copyright issues that may arise out of, say, non- or mis-attribution and other forms of plagiarism. The possibility of plagiarism through generative AI is contested by some observers, largely on the basis of how LLMs are constructed and how the information in them is made available through unique user 'prompting'. This aspect is currently unresolved in copyright law, but, as discussed above in detail, a number of lawsuits are working their way, mainly through US courts, and their outcomes should be watched closely. Users should not rule out the prospect of findings of potential copyright infringement which could impact on their work. Accordingly, in the academic research environment, the primacy of citing and referencing should be driven home, not only as good general practice, but to obviate future risk.

4.4. AI and publishing of research

Where authors have their research output published, whether as a thesis, monograph or article, they may want to take AI considerations on board. Even where the objective of publishing is more for impact than commercial return, authors should consider the consequences of commercial or OA publisher terms. More specifically, authors may have a view on whether they want their output made available for LLM training, particularly as this is a potential new source of revenue for commercial publishers. Even where publishing with open access or other freely-accessible platforms, authors should consider making some kind of rights statement, whether granting an explicit permissive waiver, such as Creative Commons CC0 1.0 Universal,¹⁰⁰ or applying 'attribution', 'commercial' or 'adaptation' constraints. Where constraints are preferred, their articulation would formally place limits on AI training whether on non-commercial or commercial bases.

library guidance. Moreover, library guides not linked here may address AI in general guidance on, e.g., research, learning, or academic integrity.

¹⁰⁰ Creative Commons. CC0 1.0 Universal Deed. See: creativecommons.org/publicdomain/zero/1.0/ (accessed 14 April 2025).

Authors should also be aware of EU rules on reserved rights under the AI Act. If they do not wish their works to be copied by AI providers, as matters stand they will have to actively opt out of the exception granted for automated computational analysis. Their publishers may or may not have covered this aspect. However, as we saw above, there is ambiguity in the relationship between the Act concerning copyright compliance and the DSM Directive concerning text and data mining. The terms of a reservation of rights must be appropriate, but this is as yet insufficiently defined, and this is enough to raise doubts on the fairness of any opt-out formulation. Specific clarification from the European Commission on this may be required.¹⁰¹

5. Conclusions

Enormous resources continue to be invested in AI. Leaving aside the latest progression from predictive chatbots to reasoning models, current AI prospects include improved potential to interrogate specialist datasets. As observed earlier, STEM and legal datasets are being curated and made available for open AI model training with technical, professional, economic and social benefits in mind. So, apart from personal and entertainment applications which drive consumer use, AI is likely to enhance capacity in science, education, commerce and government. This trend will be accelerated by the production of more stable LLMs based on reliable data sets, possibly controlled or licenced by the original content producers.

When it comes to the AI models themselves there are still more questions than conclusions. Are there limits to the capacity of predictive and especially reasoning models to generate reliable and productive output, relatively free from ‘hallucination’ and fake citations? From a research perspective, to what extent can even a properly functioning AGI improve the productivity of existing knowledge and information bases, and add to them? Is there a point at which the level of financial and computing investment fails to justify returns? Are there natural limits on the amount of verified, applicable knowledge that human-centred AI can produce and which humans can control?

The most open question from the higher education perspective is how AI will impact learning and research. Over recent decades, the online availability of full texts, data and multimedia materials, and the wide capture of search tools has greatly enhanced the ability of researchers to interrogate sources. AI has the potential to further expand the accessibility of usable content and, beyond that, to summarise sources

¹⁰¹ Mark Hyland. ‘Clash of the titans: Copyright compliance in the context of the EU’s new AI Act’. *Law Society Gazette*, December 2024, 24-29. See: <https://www.lawsociety.ie/globalassets/documents/gazette/gazette-pdfs/gazette-2024/december-2024-gazette.pdf> (accessed 14 April 2025).

and generate new content. The mechanics of a level of AI short of singularity may end up placing value on those who can interact with AI to generate outcomes, rather than those who have specialist knowledge of a subject corpus. Even today, AI competence is being marketed as a skill set which enables people who deploy it to replace not jobs per se, but those who cannot use AI.

In a positive sense, this may result in a step change in the way research is performed, e.g. by deploying specialist trained teams to interrogate LLMs and generate outputs at a faster rate. However, a trend of this nature could also undermine more traditional apprenticeship models which value post-education on-the-job learning, potentially precipitating a 'collapse of the talent pipeline'.¹⁰²

Finally, we return to the idea of copyright as property. There will always be questions around the true originality of works. However, this is built into the copyright system which, unlike patents, accepts the status of the work produced as an expression without necessarily having to be particularly original. In practice, much copyrighted work becomes a commodity with a short shelf life once published, news content and even academic research being cases in point. In this context, copyright has an economic function protecting the exploitation of works, even where the income generated, or the time span for its accumulation, may be limited. Indeed, much of the value of the protection afforded in both these examples resides in the act of attribution or citation in the context of fair use. As currently modelled, the AI process threatens these protections.

The manner of the production of LLMs has directly challenged the conventional privileges of right holders. In many instances, content, even behind paywalls, has been scraped and gathered as a fungible resource, justified on grounds of innovation. This conflict of interests may ultimately be resolved by compromise, such as compensation for use, possibly trumped by the individual right to withhold own works. This would still likely result in a scenario where AI has enough content to operate effectively at the scale it requires. However, even here, improved AI models, particularly reasoning AGIs, are likely to present a new property paradigm whereby prompting by humans, or even by autonomous AI systems themselves, will produce 'works' to which rights may theoretically accrue. In the latter case, this has the potential to conflict with the existing convention that rights accrue to a person.

Autonomous AI generation 'in the style of' a given creative person presents the additional scenario of directly undermining the ability of such human persons both to produce ongoing content and protect their existing and future economic interests. Spotify is already reported as preparing to launch an AI-powered music tool which it

¹⁰² Quote attributed to Ethan Mollick in a post on X dated 12 September 2024, which refers to passages in pages 178-180 of his book cited at footnote 9. See: <https://x.com/emollick/status/1834228159190806552> (accessed 14 April 2025).

is speculated could enable users to create personalised versions or remixes of songs.¹⁰³

For other outputs that might be characterised as less creative, and which might have a shorter exploitation life span, such as news, opinion pieces and academic articles, subordinating such content to the needs of AI may have the additional effect of commodifying them effectively at the point of publication. Without adequate safeguards, this puts AI in a position to be able to directly compete in effect immediately with research, reportage and ‘think pieces’. And this is before we even contemplate the consequences of large AI platforms concentrating output which they can market and control access to, and for which they can argue to have copyright law revised in their interests.

¹⁰³ Bernard Marr. ‘Spotify’s bold AI gamble could disrupt the entire music industry’. *Forbes*, 4 March 2025. Spotify itself claims that it does not intend to create and release music with generative technology: Emma Wilkes. ‘Spotify co-president says AI-generated music is welcome on the platform - but it won’t generate music itself’. *MusicTech*, 20 November 2024. See: <https://musictech.com/news/industry/spotify-will-host-ai-generated-music/> (accessed 14 April 2025).