

Intellectual Property Rights at the Training, Development and Generation Stages of Large Language Models

Christin Kirchhübel¹, Georgina Brown²

Atlantic Chambers, Liverpool, UK¹

Department of Linguistics and English Language, Lancaster University, UK²

ck@atlanticchambers.co.uk¹, g.brown5@lancaster.ac.uk²

Abstract

Large Language Models (LLMs) prompt new questions around Intellectual Property (IP): what is the IP status of the datasets used to train LLMs, the resulting LLMs themselves, and their outputs? The training needs of LLMs may be at odds with current copyright law, and there are active conversations around the ownership of their outputs. A report published by the House of Lords Committee following its inquiry into LLMs and generative AI criticises, among other things, the lack of government guidance, and stresses the need for clarity (through legislation, where appropriate) in this sphere. This paper considers the little guidance and caselaw there is involving AI more broadly to allow us to anticipate legal cases and arguments involving LLMs. Given the pre-emptive nature of this paper, it is not possible to provide comprehensive answers to these questions, but we hope to equip language technology communities with a more informed understanding of the current position with respect to UK copyright and patent law.

Keywords: Intellectual Property, copyright, Large Language Models (LLMs)

1. Introduction

Intellectual Property (IP) can be protected by patents, trademarks, copyright, and design rights, amongst others. As relatively uncharted territory in law, we consider copyright and patent law specifically in relation to the training and development of Large Language Models (LLMs), as well as the IP status of the outputs that LLMs generate. Because little to no IP caselaw yet exists specifically in relation to LLMs, this paper turns to discussions and caselaw concerning neighbouring Artificial Intelligence (AI) technologies as these will likely extend to LLMs as legal cases arise in the future.

2. Copyright

Copyright is an unregistered right meaning it arises automatically. It is available for literary, dramatic, musical or artistic works, sound recordings, films, broadcasts, and the typographical arrangement of published works provided that these works are 'original'. Owners of copyright have the exclusive right to do the 'acts restricted by the copyright' specified in Section 16(1) of the Copyright, Designs and Patents Act 1988 (the "1988 Act"), which includes, among others, making a copy of the work. In the absence of any defences or exceptions, copyright infringement would occur when the whole or a substantial part of copyright protected work is copied without permission, for example. When thinking about copyright in the context of LLMs, it is logical to differentiate between 'input', i.e., the

data used to train a LLM, and 'output', i.e., the data generated by a LLM.

2.1 Input

A pertinent issue is the extent to which training a LLM poses copyright infringement risks. Training a LLM relies on text and data mining (TDM) of large amounts of data. While some LLM developers are more transparent than others about the data that they rely on, there is strong evidence to suggest that, in many instances, the data used will be covered by copyright protection. For example, in written evidence to the House of Lords Communications and Digital Committee (the "House of Lords Committee") who conducted an inquiry into 'Large Language Models and Generative AI' (report dated 2 February 2024), Open AI admitted that it was "impossible to train today's leading AI models without using copyrighted materials" and attempting to do so "would not provide AI systems that meet the needs of today's citizens" (Open AI—written evidence (LLM0113)). Another example may be seen in the case of Getty Images (US), Inc. v. Stability AI, Inc., 1:23-cv-00135. Getty has issued copyright infringement proceedings (among others) against Stability AI for 'scraping' millions of images from the Getty Images Websites without Getty's consent and then using those images as input to train and develop its AI model. Getty claims that, in many cases, the output delivered by Stability AI includes a modified version of a Getty Images watermark, from which it can be inferred that Stability AI has been trained on Getty's data.

TDM is not a uniform process; rather, it varies between LLM developers. While some TDM methods may involve the copying of whole works, other TDM approaches may 'only' collect links to websites. Some may require a copy of the work to be retained, others may only necessitate temporary copies which are discarded once the relevant information has been extracted. As such, whether TDM involves copying of the whole or substantial part of the work is a moot point and likely to be case-specific.

On the assumption that TDM is considered to involve copying, it follows that developers run the risk of copyright infringement. One way for developers to avoid this risk of copyright infringement would be to rely on the exception afforded by Section 29A of the 1988 Act which permits TDM for non-commercial purposes. In the EU context, Article 3(1) of the Directive (EU) 2019/790 on copyright and related rights in the Digital Single Market appears to provide a similar exception, whereby TDM is permitted for purposes of scientific research. Needless to say, while provisions such as these may be a solution for developers operating in the research environment, they would be inadequate for commercial purposes. Another way of avoiding copyright infringement would be to obtain permission from the copyright holders by way of a licence. As already highlighted, during their inquiry into LLMs in Autumn 2023, the House of Lords Committee received oral and written evidence around the issue of copyright, and it became abundantly clear that LLM developers were using copyrighted data to train models without permission, i.e., without a licence and for commercial purposes.

In view of the evidence presented, there is a clear tension between the interests of developers on the one hand, and copyright holders on the other. There appears to be agreement between developers that access to copyright protected works is essential to ensure that AI systems perform the best they can, and the need for licence agreements could be prohibitive for this quest. However, using copyright protected works without permission goes against the whole purpose of copyright which is to reward original creations and incentivise innovation.

The UK government attempted to resolve this tension through the introduction of an AI copyright code of practice. In summer 2023, the UK Intellectual Property Office (UKIPO) set up a working group involving stakeholders from the technology, creative and research sectors. Members included the BBC, the British Library, Financial Times, Google Deepmind, IBM, Microsoft, Stability AI, UK Research and Innovation. The UKIPO said that: *'The code of practice aims to make licences for data mining more available. It will help to overcome barriers*

that AI firms and users currently face, and ensure there are protections for rights holders. This ensures that the UK copyright framework promotes and rewards investment in creativity. It also supports the ambition for the UK to be a world leader in research and AI innovation.' (UKIPO, 2023). In February 2024, the UK government announced that it had shelved plans to put in place this code as it had become clear that the working group would not be able to reach agreement. (Department for Science, Innovation & Technology, 2024: 19).

While we wait for the government to clarify the relationship between IP and AI, all we are left with is the existing law. Considering this issue in the context of the 1988 Act, there is legal uncertainty as to whether commercial AI developers are infringing copyright when training AI systems on copyright protected material without a licence. As alluded to above, it may be argued that TDM does not actually involve 'making a copy' for the purposes of the 1988 Act, and even if it did, this copy may only be temporary and therefore fall within the exceptions allowing for transient or incidental copies provided for by Section 28A of the 1988 Act. Further, even if TDM were to involve copying, given that this is only part of the training stage, and given that the actual AI model does not directly reproduce the copyright protected work, but rather reflects the data / information contained within that work, does copyright even apply in those circumstances? Understandably, litigation around issues such as these has already started and is likely to grow in the near future.

2.2 Output

Questions around copyrightability also arise for LLM output – can AI-generated material attract copyright protection in the first place, and if so, who is the author? In relation to the former, the 1988 Act expressly provides for the copyright protection of literary, dramatic, musical or artistic work which is computer-generated in circumstances such that there is no human author of the work. This is in contrast to the position in the US, for example, where only works with human authors can receive copyright protection. Even though current UK legislation seems to explicitly cater for copyright in AI-generated work, in order for this work to be a true candidate for copyright protection, it needs to be 'original'.

Traditionally, the test for originality in the UK had a low threshold, requiring the work to be produced with 'sufficient skill, labour, and judgment'. This changed following the judgment of the European Court of Justice in the case of *Infopaq International A/S v Danske Dagblades Forening* [2009], with UK courts adopting the EU requirement of work having to exhibit the 'author's own intellectual creation' in order to be deemed

'original'. There is some uncertainty around how UK courts are going to interpret the originality requirement going forward. However, in view of provision 9(3) of the 1988 Act, which states, '*In the case of a literary, dramatic, musical or artistic work which is computer-generated, the author shall be taken to be the person by whom the arrangements necessary for the creation of the work are undertaken*', it would appear that works without a human author could meet the originality requirement. It follows that AI-generated work which has had some element of human involvement may very well pass the originality test and therefore could be copyrightable.

Turning now to the question of authorship, there are a number of possible authors – the developer of the AI system, the user of the AI system, i.e., the prompt engineer, or the AI system itself. Section 9(3) of the 1988 Act (reproduced above) provides that the author of computer-generated work '*shall be taken to be the person by whom the arrangements necessary for the creation of the work are undertaken*.' Given that statutory drafting refers to a '*person*', it may be the case that the AI system cannot be the author for the purposes of the 1988 Act. In the absence of a contract, whether it is the person who built the AI system, or the prompt engineer 'who made the necessary arrangements' is likely to be case-specific.

As with the input stage, the topic of copyright infringement is also relevant to the output stage. Guadamuz (2024) presents a detailed discussion of the relevant issues including possible defences to arguments that AI generated output infringes copyright. Briefly here, the points that are likely going to be debated in this sphere include the copyright infringement potential of memorisation, i.e., LLM models "memorising" specific fragments of their training data, and then reproducing these fragments in their output (Emanuilov and Margoni, 2024). Rather than reproducing existing work 'verbatim', perhaps a more likely scenario involves AI output resembling input data so the legal analysis will revolve around similarity of input vs. output.

We do not have existing legal authority on these issues, but this is likely to change in the near future in view of *Getty Images v Stability AI* [2023] EWHC 3090, a claim which is currently in the process of being litigated in the High Court. The Getty case will be of particular interest as it raises IP right infringement issues around input as well as output.

3. Patents

There is a question of whether the outputs of AI can be patented. Patents fall within the so-called registered rights as they are granted on application to the UKIPO. Patents provide the patent holder with an exclusive right over the

invention, e.g., an exclusive right to a product or a process, for a period of time. In order to qualify for a patent, the product or process must be new, involve an inventive step, be capable of industrial application, and not specifically excluded from protection. In exchange for the patent grant, the applicant must disclose technical information about the invention to the public.

While issues around AI and copyright remain to be tested in the courts, we do have some legal authority in the area of patent law.

3.1 Can AI be an inventor?

The question of whether AI can be an inventor under current UK patent law has already been considered by the Supreme Court in the case of *Thaler v Comptroller-General of Patents, Designs and Trade Marks* [2023] UKSC 49. Dr Thaler filed two patent applications under the Patents Act 1977 (the "1977 Act") for inventions solely created by an AI system called DABUS of which Dr Thaler was the owner. The Hearing Officer for the Comptroller-General of Patents, Designs and Trade Marks (the "Comptroller") issued a decision that (i) DABUS could not be an inventor for the purposes of Sections 7 and 13 of the 1977 Act because it was not a person, and (ii) Dr Thaler was not entitled to a patent based on his ownership of DABUS in circumstances where DABUS was listed as the inventor. Dr Thaler appealed to the High Court and the Court of Appeal but both appeals were unsuccessful.

The Supreme Court decided that: i) an inventor within the meaning of the 1977 Act must be a natural person, and ii) that the doctrine of accession does not apply as this is not a case where new tangible property is produced by an existing item of tangible property. It follows that DABUS is not an inventor for the purposes of 1977 Act and the Act did not confer on Dr Thaler the property in or the right to apply for and obtain a patent for any technical development made by DABUS. Accordingly, the Comptroller was right to find Dr Thaler's applications as withdrawn under Section 13(2) of the Patents Act. The Supreme Court acknowledged that had it been Dr Thaler's case that he was the inventor (rather than DABUS), and that he had used DABUS as a '*highly sophisticated tool*', the outcome of the proceedings '*might well have been different*'.

It is important to note that, right at the outset of the judgment, Lord Kitchin (with whom the other Lords agreed) made clear that '*the appeal is not concerned with the broader question whether technical advances generated by machines acting autonomously and powered by AI should be patentable. Nor is it concerned with the question whether the meaning of the term "inventor" ought to be expanded, so far as necessary, to include*

machines powered by AI which generate new and non-obvious products and processes [48] ... This appeal is concerned instead with the much more focused question of the correct interpretation and application of the relevant provisions of the 1977 Act to the applications made by Dr Thaler.' [50] The court recognised that, in view of rapid advances in AI technology, these broader questions are increasingly important and alluded to a potential shift in the legal landscape as a result. However, in citing the judgment of Laing LJ in the Court of Appeal, the Supreme Court was clear that *'if patents are to be granted in respect of inventions made by machines, the 1977 Act will have to be amended'* [79].

Whether AI can be an inventor for the purposes of patent law has received global attention with Dr Thaler filing test patent applications in different jurisdictions around the world, including the European Patent Office (EPO). The current legal position on an international level appears to be that an inventor for patentable inventions must be a human or a person with legal capacity.

3.2 Can AI be patented?

The inventions in the Thaler case concerned a food container and a light beacon. It was undisputed that these were patentable, i.e., there was no issue around novelty, for example, and the inventions did not fall within categories that are excluded from patentability. What about the patentability of the AI system itself? Under Section 1(2)(c) of the Patents Act 1977 'a programme for a computer ... as such' is excluded from patent protection. Essentially, the position is that one cannot obtain a patent for a computer programme in itself; however, if the computer programme provides a 'technical contribution' to the real world, then it is patentable. A recent decision of the High Court considered this statutory provision and associated caselaw in the context of AI in the case of Emotional Perception AI Ltd v Comptroller-General of Patents, Designs and Trade Marks [2023] EWHC 2948.

Emotional Perception had applied to patent an Artificial Neural Network (ANN). The ANN was said to be capable of providing improved media file recommendations. Taking the context of music websites for example, where a user wants to receive music similar to music they already have, traditional tools would recommend music tracks based on similar categories of music (e.g., rock), those categories having been tagged as such by humans. Instead of taking the 'category' of music as the criterion for recommending similar music tracks, the ANN-based system is said to identify similar music tracks based on human perception and emotion. In brief, the system works as follows: it takes a pair of music files,

which have been given a semantic label, e.g., 'happy', 'relaxing', and so on. The files are plotted in a 'semantic space', with the distance between the files indicating their semantic similarity. In addition to their semantic properties, the files are also analysed according to physical properties such as tone, timbre, speed etc., and again plotted in a 'property space'. Using back-propagation, the property space is refined in order to reflect the semantic space, so that semantically similar tracks are close in property space, whereas semantically dissimilar tracks are farther apart in the property space. The operational ANN is then able to take a music track, determine its physical attributes, plot these against the physical attributes of other music tracks in a music library or database, and by looking for those tracks which are most proximate in terms of physical characteristics, it can recommend semantically similar tracks.

An officer for the UKIPO refused grant of the patent on the basis that the ANN system was considered to be 'a program for a computer' and that the patent application related to that computer programme 'as such'.

Emotional Perception appealed to the High Court challenging the decision by the UKIPO to refuse grant of the patent. The matter came before Sir Anthony Mann, J, who considered whether i) the ANN was 'a program for a computer' therefore falling within the statutory exclusions to patentability, and ii) if it was, whether there was a technical contribution which meant it fell outside the exclusionary regime.

On the first point, Mann J, differentiated between hardware ANNs and software emulated ANNs, and concluded that neither qualifies as 'a programme for a computer' and therefore neither was excluded from patentability. In the case of hardware ANNs, it was accepted by the parties that there is no 'programme' and therefore this would not fall within the exclusions. *'The hardware is not implementing a series of instructions pre-ordained by a human. It is operating according to something that it has learned itself.'* [54] In the case of software emulated ANNs, there were two aspects in which computer programming plays a role, one being the training stage, and the other being the software platform which enabled the computer to carry out the emulation. With regards to the latter, Mann, J considered that this can be de-coupled from the ANN: *'It seems to me that it is appropriate to look at the emulated ANN as, in substance, operating at a different level (albeit metaphorically) from the underlying software on the computer, and it is operating in the same way as the hardware ANN. If the latter is not operating a program then neither is the emulation.'* [56]

The court found that the ANN, in itself, was not a computer programme because it was not

operating a set of programme instructions given to it by a human. The ANN had trained itself, applying its own weights and biases. It was emulating a piece of hardware which had physical nodes and layers, and was no more operating or applying a program than a hardware system was.

With respect to the computer programme involved at the training stage, which sets the training objectives and parameters in which the ANN is to operate, the court concluded that this fell outside the actual invention that is claimed. The invention was not a claim to the computer programme at the training stage; the invention related to the idea of using pairs of files for training, and setting the training objective and parameters accordingly. The claim therefore went beyond the actual computer programme.

As explained above, even if an invention were to be a claim to a computer program, it may still be patentable if it provides a “technical contribution” outside the computer program itself. Given his conclusion that the ANN was not a computer programme, Mann J, did not need to consider the question of technical contribution, but he nevertheless did. Following a review of caselaw on what constitutes a ‘technical contribution’, Mann J found that the sending of a file recommendation to an end user is a matter external to the computer and amounts to a technical contribution, i.e., the ANN has a real world effect outside of the computer.

So far, established practice at the UKIPO was to treat inventions involving AI as computer implemented and therefore applications would have had to be considered under the computer program exclusion exemption, i.e., whether the invention produced a technical contribution. The position in Europe has been similar with the EPO considering inventions involving AI as computer-implemented inventions which would only become patentable if they are applied to solve a technical problem in a field of technology.

The Emotional Perception AI judgment has potentially opened up a new avenue to obtain patent protection in the UK for inventions involving ANNs and AI more generally. We say ‘potentially’ as UKIPO is currently appealing the decision of the High Court. We will await to see whether the Court of Appeal, like the High Court, reaches a decision favourable to patentees of AI inventions. In case the High Court decision is upheld, it will be interesting to see what influence the Emotion Perception AI judgment will have on the approach taken by the EPO.

4. Conclusion

The training of AI technology has led to copyright disputes, and there are question marks over the IP of the outputs that result from generative AI (and the IP status of the AI itself). It is quite easy

to see how the cases and discussions drawn on in this paper extend to LLMs. While we await further development and resolution in these cases and discussions, this paper has aimed to put a spotlight on the issues that could feasibly arise for LLM stakeholders going forward.

5. References

5.1 Academic

Emanuilov, I., & Margoni, T. (2024). Forget me not: memorisation in generative sequence models trained on open source licensed code. DOI: 10.5281/zenodo.10635479

Guadamuz, A. (2024). A Scanner Darkly: Copyright liability and exceptions in Artificial Intelligence inputs and outputs. GRUR International. 73(2) pp. 111-127 DOI: 10.1093/grurint/ikad140.

5.2 Legislation and Caselaw

Copyright, Designs and Patents Act 1988

Directive (EU) 2019/790 on copyright and related rights in the Digital Single Market

Emotional Perception AI Ltd v Comptroller-General of Patents, Designs and Trade Marks [2023] EWHC 2948

Getty Images v Stability AI [2023] EWHC 3090

Getty Images (US), Inc. v. Stability AI, Inc., 1:23-cv-00135

Infopaq International A/S v Danske Dagblades Forening [2009] ECJ Case C-5/08, ECLI:EU:C:2009:465

Patents Act 1977

Thaler v Comptroller- General of Patents, Designs and Trade Marks [2023] UKSC 49

5.3 Guidance

Communications and Digital Committee Large language models and generative AI. 1st Report of Session 2023-24. HL Paper 54. 2nd February 2024. URL: <https://publications.parliament.uk/pa/ld5804/ldslect/ldcomm/54/5402.htm>

Department for Science, Innovation & Technology. Consultation Outcome: A pro-innovation approach to AI regulation: Government response. Command paper CP1019. 6th February 2024. URL: <https://www.gov.uk/government/consultations/ai-regulation-a-pro-innovation-approach-policy-proposals/outcome/a-pro-innovation-approach-to-ai-regulation-government-response>

Open AI Written Evidence to HL Paper 54. URL: <https://committees.parliament.uk/writtenevidence/126981/html/>

UK Intellectual Property Office. (2023). The governments code of practice on copyright and AI. URL: <https://www.gov.uk/guidance/the-governments-code-of-practice-on-copyright-and-ai>