

PUBLICATION

The Fast-Moving Race Between Gen-AI and Copyright Law

Authors: Scott M. Douglass, Dominic A. Rota

July 10, 2024

With the emergence of Generative AI platforms such as ChatGPT and Google Gemini, the widespread use of copyrighted works to train the software systems behind these programs is raising pressing legal questions about the permissible use of copyrightable works. Recently, there has been a spree of licensing deals between businesses that develop generative AI systems and those that own the content necessary to train those systems. OpenAI, the developer of ChatGPT, has signed deals with News Corp., Vox Media, The Atlantic, Thomson Reuters, Associated Press, and several other news publishers. It is also in talks with CNN and Time, reportedly offering \$1 million to \$5 million per year to license content to train its AI systems. While seemingly expensive, this offer is relatively low compared to the \$250 million over five years that OpenAI agreed to pay News Corp.

At the same time, content creators (including Sarah Silverman, Getty Images, and The New York Times) are suing developers of generative AI systems, most prominently but not exclusively OpenAI, for using their copyrighted works to train those systems without licensed permission. Most recently, the Center for Investigative Reporting (CIR) filed a complaint against OpenAI and Microsoft on June 27, 2024, alleging that they "copied, used, abridged, and displayed CIR's valuable content" without permission or authorization and without providing compensation. Universal Music Group, Sony Music Entertainment, and Warner Records also filed complaints on June 24, 2024, against AI companies Suno and Udio, alleging violation of their copyrights "en masse" in the development of AI music generators. These lawsuits are all fairly new, most are no further than the pleadings stage, and are ongoing, so the issues of whether generative AI training data and outputs constitute copyright infringement and whether fair use is a viable defense are still in flux.

This article analyzes the impact of emerging generative AI technology on existing copyright law, discusses pending litigation regarding each of these issues (direct copyright infringement, derivative use, fair use, and violation under the Digital Millennium Copyright Act), and analyzes the implications for copyright owners and developers of generative AI technology.

Generative AI and Copyright

Generative AI systems, such as ChatGPT, are able to generate new and different content in response to a text prompt entered by the user. These programs rely on large language models (LLMs), which are trained by inputting enormous amounts of data, including images, text, audio, and video gathered from the internet. To be trained effectively, generative AI models require vast amounts of data because the system "is only as good as its training data" ([Forbes](#)).

Currently, there is no mechanism to track and share what inputs generative AI models have ingested during the training process. This raises several legal issues, especially if some or all of the training data is subject to copyrights. Because training the models usually involves making digital copies, companies developing generative AI systems could be reproducing copyrighted works: a *prima facie* case of copyright infringement. The outputs generated by the AI programs could also be considered derivative works of original copyrighted works, constituting an indirect form of copyright infringement. Furthermore, if the LLMs remove copyright management information (CMI) from the inputted training data, the copyright holder could have a claim under the Digital Millennium Copyright Act (DMCA).

In this uncharted legal territory, advocates of AI technology and creators' rights debate the proper course of action: generative AI is either infringing on copyright owners' exclusive rights or is utilizing the copyrighted works in a way that constitutes fair use. If the former is true, AI companies will have to license all the input data that are subject to copyrights, creating practical and logistical problems that present an "existential threat" to generative AI technology ([Bloomberg](#)). On the other hand, if the use of copyrighted works to train generative AI programs is considered fair use, meaning there is no requirement to obtain a license, the creators are not fairly compensated for the use of their protected expression. Courts have only just begun to wade into this debate with the recent deluge of litigation.

Direct Infringement by Reproduction

One of the three major claims plaintiffs raise is direct copyright infringement through the reproduction, or copying, of copyrighted works to train generative AI systems. To establish a *prima facie* claim of infringement, the plaintiff must show ownership of a valid copyright and copying of the original elements of the work that are protected by copyright. A few claims of infringement against AI companies have already progressed past the motion-to-dismiss stage, but the courts' rulings shed little light on the ultimate resolution of these claims.

In a case against Stability AI, the Northern District of California denied Stability AI's motion to dismiss because the question of whether copying occurred in the training or running of its generative AI program, Stable Diffusion, could not be resolved at the motion-to-dismiss stage. In one of the copyright infringement claims, the plaintiff, Andersen, relied on the output of a search on "ihavebeentrained.com," a resource for identifying the training data used by a generative AI system to produce a particular output. The Court found this to support a reasonable and plausible inference that her copyrighted works were used to train Stable Diffusion, even if the plaintiff did not specifically identify which of the plaintiff's works were used. This ruling suggests that copyright holders may be able to pass the first hurdle of an infringement claim by simply relying on this publicly accessible tool, without needing to show specific data of which copyrighted works were used.

The class action by Sarah Silverman and other authors against OpenAI alleges that OpenAI copied the plaintiffs' copyrighted books to use in the training dataset for its LLMs. OpenAI did not move to dismiss the direct copyright infringement claim regarding the copying of books to use as training data. However, the Court dismissed the vicarious copyright infringement claim, with leave to amend, because the plaintiffs did not show the second required element of an infringement claim that ChatGPT directly copied their books in its outputs or that the outputs are substantially similar to their copyrighted works. Plaintiffs filed an amended complaint, omitting the claim for vicarious infringement. Thus, whether OpenAI is vicariously or directly committing copyright infringement, is still an open question that will very likely not be resolved by this dispute.

AI-Generated Works = Derivative Works?

The use of copyrighted works to train generative AI models can also give rise to derivative use claims because AI systems produce an output that may be similar to the input training data. The statutory definition of derivative work is "a work based upon one or more preexisting works, such as a translation, musical arrangement, dramatization...or any other form in which a work may be recast, transformed, or adapted." Because a derivative work is subject to the copyright of the original work, the copyright owner of the original work has a right of action against the producer of the derivative work who did not obtain permission. For derivative use claims, courts determine whether the secondary work is substantially similar to or represents protected aspects of the original copyrighted work. For example, in *Andersen v. Stability AI Ltd.*, the Court dismissed Andersen's derivative use claim with leave to amend because the plaintiff did not allege substantial similarity or representation of protected aspects. As another example, in Sarah Silverman's case against OpenAI, the plaintiffs were permitted to amend the complaint to include the requisite allegation to support a derivative use claim. While there is still no resolution regarding whether the output of a generative AI model can constitute a derivative work, given the factual similarities with man-made derivative works, it seems likely courts may apply the same "substantial similarity" standard to AI-generated derivative works.

Fair Use Defense

Assuming plaintiffs can show a *prima facie* claim of copyright infringement, AI platform providers may raise the fair use defense, which would absolve them of liability if successful. Fair use is the right to use a copyrighted work without consent or license from the copyright owner in certain circumstances. To determine whether an infringing work constitutes fair use, courts apply a four-factor balancing test statutorily enumerated in the Copyright Act.

The fair use doctrine has not yet been applied in the context of generative AI. In looking to guidance from prior judicial decisions, including a 2015 Second Circuit decision, it is likely that the application of the four factors could weigh against a finding of fair use. As to the first factor, generative AI platforms are generally offered for a commercial purpose. Second, generative AI systems have the capacity to produce new works that closely resemble the originals. Third, generative AI technologies are trained with the whole of the copyrighted work. And, fourth, AI has the potential to generate an effective substitute for copyrighted work in the marketplace. These factors would especially favor plaintiffs if their copyrighted works are of a creative focus, rather than factual because creative works often comprise more protectable expression.

Most of the cases currently pending have not yet discussed the fair use defense. The order denying motions for summary judgment in *Thomson Reuters Enterprise Centre GmbH et al v. Ross Intelligence, Inc.*, only briefly mentions Ross's defense. That opinion states that Ross's use of Westlaw headnotes would not be transformative if Ross used the exact text to make the AI system replicate and reproduce the protected expression of the headnotes, but it would be transformative if the program only studied the language patterns to learn how to produce a similar result. The judge also states that Ross's AI program likely satisfies the amount and a substantiality factor because it only produces the judicial opinion, not the creative expression of the headnotes. The judge left the other two factors of fair use for the jury to decide, so it remains unclear which side the factors will favor when considered in aggregate. Ultimately, we await future guidance from the courts, as the fair-use doctrine is applied in a fact-intensive manner and is thus not typically adjudicated at the pleadings stage. (See our colleagues' previous discussion of fair use implications of generative AI, [here](#).)

Claims Under the Digital Millennium Copyright Act

When training their generative AI tools with copyrighted works, AI companies may also violate the Digital Millennium Copyright Act (DMCA) by removing the copyright management information (CMI) from the inputted data. In order to bring a claim under the DMCA, copyright owners must identify what the removed or altered CMI was, and that the defendant knew, or had reasonable grounds to know, that intentionally removing CMI would induce, enable, facilitate, or conceal infringement. These requirements present difficulties for copyright owners because they have to identify the particular removal of CMI among millions of works included in training datasets.

The DMCA claim in Silverman's case was dismissed with leave to amend because the plaintiffs failed to plead factual allegations supporting their claim that OpenAI removed CMI from the copyrighted books used for training and that OpenAI knew, or had reasonable grounds to know, that ChatGPT's output would induce, enable, facilitate, or conceal infringement. Similarly, the *Andersen* Court dismissed the DMCA claim with leave to amend because the plaintiffs only alleged DMCA liability for the defendants generally when they should have identified the specific CMI that was removed during training by each defendant. Plaintiffs claiming a DMCA violation therefore have a daunting task to catalogue all their copyrighted works used in training in order to survive a motion to dismiss.

Conclusion

It is still an open question whether plaintiffs will succeed in showing that use of copyrighted works to train generative AI constitutes copyright infringement and be able to overcome the fair use defense or succeed in showing that generative AI developers are removing CMI in violation of the DMCA.

The government has made some moves in the past few months to resolve these issues. The U.S. Copyright Office started an [inquiry](#) in August 2023, seeking public comments on copyright law and policy issues raised by AI systems, and Rep. Adam Schiff (D-Calif.) introduced a new bill in April 2024, that would require people creating a training dataset for a generative AI system to submit to the Register of Copyrights a detailed summary of any copyrighted works used in training. These initiatives will most likely take some time, meaning that currently pending litigation is vitally important for defining copyright law as it applies to generative AI.

Recent licensing deals with news publishers appear to be anywhere from \$1 million to \$60 million per year, meaning that AI companies will have to pay an enormous amount to license all the copyrighted works necessary to train their generative AI models effectively. However, as potential damages in a copyright infringement case could be billions of dollars, as claimed by Getty Images and other plaintiffs, developers of generative AI programs should seriously consider licensing any copyrighted works used as training data.

If you have questions about content copyright licensing or would like assistance reviewing your intellectual property portfolio or policy for AI governance, reach out to [Scott M. Douglass](#), [Dominic Rota](#), or any member of [Baker Donelson's Intellectual Property Team](#).

Charlotte Brownell, a summer associate at Baker Donelson, contributed to this alert.