

Generative Artificial Intelligence, LLMs, And Fair Use After *Warhol*: The Copyright Office and Accountability

By Thomas L. Hamlin*

Introduction

Generative Artificial Intelligence (AI)¹ has been characterized as the greatest invention since the Internet, the new new thing, a plagiarism engine, and a technology that will destroy civilization. Even Sam Altman, the founder and CEO of OpenAI, has called it an “alien intelligence.” While some or all of these descriptions may or may not be accurate, one thing is abundantly clear: the technology raises serious copyright infringement and fair-use issues that the United States Copyright Office must address to introduce accountability to a handful of Silicon Valley entrepreneurs and large tech companies who control this technology. The Office has begun this process by recently publishing “Registration Guidance: Works Containing Material Generated by Artificial Intelligence,” which presented criteria for granting copyrights to works generated by AI, as long as the product was the result of human authorship.

The Office has also held Listening Sessions this spring in which representatives from the literary, music, software, and visual-art worlds offered opinions about how the Copyright Office should address this new technology while protecting the rights of creatives in these various industries. The Listening Session participants offered a variety of comments, but the major concern was that AI companies “scrape” the internet for huge volumes of information to train Large Language Models (LLMs), which, in turn, power chatbots such as ChatGPT. Much of this information is protected by copyrights, but these AI companies offer no compensation to creators. This raises the issue of whether training these LLMs and producing their outputs infringe copyrights under 17 U.S.C. § 601, or are “fair use” under 17 U.S.C. § 701. The Copyright Office has not yet proposed regulations on the training of LLMs, or when fair use may apply. It plans to do that after seeking even more comments from the various stakeholders, and other interested parties. While the Office itself does not litigate fair-use issues, it does publish a best-practices guideline on fair use, as part of its statutory mandate to administer the Copyright Act.

Recently, the United States Supreme Court issued a decision on at least one of the statutory factors for fair use that must be considered as a defense to copyright infringement: “the purpose and character of the use, including whether such use is of a commercial nature or is for nonprofit educational purposes.”² While the facts in *Andy*

* Mr. Hamlin is a senior partner with Robins Kaplan LLP in its Minneapolis office. He has received State and National awards and recognition for his work in the areas of intellectual property and business litigation.

Warhol Found. for the Visual Arts v. Goldsmith did not involve AI, the Court’s extensive discussion of the purpose and contours of the fair-use doctrine, particularly that commercial use may be more important than the transformative nature of the use in determining whether fair use applies, nevertheless has relevance to the issue of whether LLMs and their output are entitled to fair-use protection.

This article will begin with a definition of AI provided by the technologists who created it. Then it will turn to a discussion of the Copyright Office’s recent guidance on AI-generated works, and a summary of comments and opinions provided to the Copyright Office from industry representatives on training LLMs and fair use. The article will discuss the Supreme Court’s analysis of fair use in the *Warhol* decision, and its potential impact on training LLMs on copyrighted works, including producing outputs based on those works. Recent cases applying *Warhol* will also be discussed, along with two cases relied on by those who claim training LLMs is fair use. The article will conclude with a proposal to train LLMs while protecting copyrights – a proposal that the Copyright Office can use to add further protections to creative works, and to impose accountability on a nascent but increasingly important industry.

What Exactly Is Generative Artificial Intelligence?

Generative artificial intelligence has been described as the use of machines to mimic human intelligence. At this point in its development, however, this technology is not a substitute for the human brain. Instead, it is a multitude of neural networks that absorb huge sums of data in parallel. It then codifies nearly every pattern of human language through the use of complex algorithms, and generates content in text, images, code, video, or audio. No human brain has the capacity to do what AI does. While AI can absorb the entire history of the human word, that is far beyond human capabilities. At the same time, AI is capable of “learning” through a trial and error process of prediction.

OpenAI’s chief scientist, Ilya Sutskever, has explained that “[a] neural network learns, and its learning is powered by prediction – a bit like the scientific method. The neurons sit in layers. An input layer receives a chunk of data, a bit of text or an image, for example. The magic happens in the middle – or ‘hidden’ – layers, which process the chunk of data, so that the output layer can spit out its prediction.... A neural network learns because its training data include the correct predictions, which means it can grade its own outputs.... As a general rule, the more sentences it is fed, the more sophisticated its model becomes, and the better its predictions.”³

Mr. Sutskever’s description of AI underscores the importance of training LLMs using reams of information – the more information, the better the “prediction” or output. What’s more, the enormous volume of information needed by LLMs comes from software that “scrapes” or gathers from the internet this information, much of it protected by copyright, that forms the basis of the AI output. From this information, AI

learns to recognize patterns, structures, and context in human language, which LLMs can then use to emulate written text. This raises the critical issue of whether the compilation and use of copyrighted works to train LLMs and produce its output is infringement or fair use – an issue the Copyright Office is currently grappling with, and one that is also the subject of ongoing litigation.

The Copyright Office is not the only branch of the federal government focusing on AI technology. On October 30, 2023, President Biden issued an Executive Order on the Safe and Secure, and Trustworthy Development and Use of Artificial Intelligence, which, among many other directives, instructed the Secretary of Commerce to submit a report on the “potential benefits, risks, and implications” of the use of LLMs.⁴ Thus, the Copyright Office is part of a multipronged effort to place safeguards around this revolutionary technology. This article, however, will limit its discussion to AI issues confronting the Copyright Office.

The Copyright Office Has Issued Guidance on Works Containing Material Generated by AI.

On March 16, 2023, the Copyright Office issued a “Registration Guidance: Works Containing Material Generated by Artificial Intelligence.”⁵ The Office stated that “applicants have a duty to disclose the inclusion of AI-generated content in a work submitted for registration and to provide a brief explanation of the human author’s contribution to the work. As contemplated by the Copyright Act, such disclosures are ‘information regarded by the Register of Copyrights as bearing upon the preparation or identification of the work or the existence, ownership, or duration of the copyright.’” The Office emphasized that “[i]n each case, what matters is the extent to which the human had creative control over the work’s expression and ‘actually formed’ the traditional elements of authorship.”

The Registration Guidance goes on to explain that “to qualify as a work of ‘authorship’ a work must be created by a human being, and that it will not register works produced by a machine or mere mechanical process that operates randomly or automatically without any creative input or intervention from a human author.... In the case of works containing AI-generated material, the Office will consider whether the AI contributions are the result of ‘mechanical reproduction’ or instead of an author’s ‘own original mental conception,’ in which the author gives visible form. The answer will depend on the circumstances, particularly how the AI tool operates and how it was used to create the final work. This is necessarily a case-by-case inquiry.”⁶

The Registration Guidance did not offer comments on the training of LLMs: “the Office intends to publish a notice of inquiry later this year seeking public input on additional legal and policy topics, including how the law should apply to the use of copyrighted works in AI training and the resulting treatment of outputs.”⁷

The Copyright Office's Listening Sessions On the Use of AI

The Copyright Office held a series of “Listening Sessions” in the spring of 2023 to gather input on AI and copyrights from stakeholders in the diverse worlds of literature, software, music, and the visual arts. This article will focus on only one of the Listening Sessions, the one exploring comments on literary works, including software, because the participants’ comments reflect many of the concerns expressed in other Listening Sessions. The group in the literary works Listening Session included writers, publishers, and software developers, as well an academic and an attorney representing a prominent hedge fund that invests in AI companies. This session was held on April 19, 2023, before the Supreme Court issued the *Warhol* decision. One of the questions the Copyright Office asked the participants to comment on was: “How is training or output of AI affecting your industry?”⁸ The responses varied depending on the industry representative. Below is a sample of the comments expressed in the literary works session.

- A member of the Copyright Alliance, which represents individual creators and organizations across the spectrum of copyright disciplines observed: “[T]he Copyright Alliance supports the responsible and ethical advancement of AI technology. Many in the creative industry are already using or planning to use AI for the creation of a wide range of works that benefit society.... [But] [t]he interests of those using copyright materials must not be prioritized over the rights and interests of creators and copyright owners.... [B]oth small and large creators face significant risk of being harmed when their works are copied without their authority for ingestion purposes. In particular, individual creators, who have little to no negotiating power with AI system developers, are most at risk of such harms.”⁹
- A spokesman for the National Writers Union, which includes writers in all genres and media noted: “Our members have created works which have been scraped from the internet, copied, and used for training generative AI without permission or payment....Congress could best facilitate organizing, collective bargaining, and collective licensing for AI training by explicitly clarifying the right of freelancers and self-publishers to organize and act collectively as workers, including but not limited to collective bargaining over the terms of collective licensing.”¹⁰
- An Associate Professor of Computer and Information Science offered a contrary opinion: “[R]egarding guidance on whether training AI systems on copyright materials without affirmative consent from the right holders should be considered fair use, I’d like to argue that it should be considered fair use because,

first, fair use or learning, if the Copyright Office decided that it were not fair use, then that would make training of these AI systems effectively impossible and would shut down this interesting development. Second, the learning process is transformative.”¹¹

- The CEO of the Authors Guild advocated for its members: “[S]o GPT and Bard, the two main engines, were developed by copying and ingesting large amounts of texts, including potentially millions of books and articles found online without permission, and in generating the outputs, these programs merely re-scramble inputs. Nothing new is added. Generative AI cannot think or feel itself. It cannot express emotion. It can only mimic what it has been fed. And so, by its nature it is always derivative of what it’s been trained on. There would be no GPT without pre-existing works...We believe that this use is not fair and that it should be compensated. We do not, however, want to impede the development of AI, so we would like to see collective licensing that makes it possible for AI developers to license the data they need and compensate authors.”¹²
- An attorney for a large venture capital fund that invests in companies that both build and rely on AI offered his client’s viewpoint: “So the point I want to be sure to emphasize is that, really the only practical way for these tools to exist is if they can be trained on massive amounts of data without having to license that data. In fact, the data needed is so massive that even collective licensing really can’t work. What we’re talking about in the context of these large language models is training a corpus that is essentially the entire volume of the written word.”¹³
- A spokesperson for the News Media Alliance, which represents the most trusted publishers in print and digital media defended the rights of its members: “These systems have been developed by ingesting massive amounts of the creative output of publishers, often without authorization or compensation, and they disseminate that same content in response to user queries, again without authorization or payment and often with little or no attribution or link to the original news source....Copyright laws should protect and not harm publishers in this setting. Developers and deployers of generative AI should not use expressive works without authorization and should respect publishers’ rights to negotiate fair compensation for the use of their valuable works. The system should also be transparent to publishers and users. They should identify the content used to fuel their products and connect and not disintermediate users with publishers. Protecting publishers’ legitimate intellectual property interests will strengthen, not impede, generative AI innovation because authorized use of publisher content can improve the reliability and accuracy of AI products, which will enhance system output and bolster consumer confidence.”¹⁴

- The attorney for the private equity fund added that imposing costs on AI creators will have several negative effects: “[I]f we’re thinking about imposing new costs on creators of AI models, I think one of two things is going to happen. I think either these tools just won’t be able to be built, and I think that’s probably the most likely outcome because, because (sic) of the way these tools are built, they require just way too much data for any licensing scheme to be able to work. At best, what will happen is that the ability to build these tools will be preserved for those companies that have the deepest pockets and the greatest incentive to keep AI models closed, so the result will be less competition, far less innovation, and closed AI models, which are hard to investigate. So I think we ought to be very, very cautious about imposing new costs on the creators of these new tools....”¹⁵

Andy Warhol Found. for the Visual Arts v. Goldsmith

Lynn Goldsmith, an accomplished, professional photographer, was commissioned by Newsweek in 1981 to photograph an up-and-coming musician named Prince for use in an article about him. Ms. Goldsmith’s black and white portrait of the artist was the copyrighted work at issue in the court case. Several years later, Ms. Goldsmith granted a limited license to Vanity Fair for use of one of her photographs of Prince. The terms included the phrase, “for one time only.” Vanity Fair then commissioned the renowned artist, Andy Warhol, to create a purple silkscreen portrait of the musician which appeared in its magazine. In addition to the Vanity Fair-licensed photograph, Warhol created 15 other works based on the photograph, which are collectively referred to as the “Prince Series.” The Series passed to the Warhol Foundation after Mr. Warhol died. After Prince died in 2016, Conde Nast, Vanity Fair’s owner, purchased a license from the Warhol Foundation to publish the Orange Prince for the magazine, without knowledge or consent from Ms. Goldsmith.¹⁶ The *Warhol* Court emphasized “that Goldsmith, too, had licensed her Prince to magazines such as Newsweek....,” including magazines such as People and Rolling Stone that did tributes to Prince after his death. With the exception of Conde Nast, those magazines credited Ms. Goldsmith’s photograph.¹⁷

The Court provided a simple description of the Conde Nast photograph: “Orange Prince crops, flattens, traces, and colors the photo but otherwise does not alter it.”¹⁸ When Ms. Goldsmith saw the new photograph, she notified the Warhol Foundation that it had infringed her copyright. The Warhol Foundation sued Ms. Goldsmith seeking a declaratory judgment of no copyright infringement and fair use, and Ms. Goldsmith countersued alleging infringement. The District Court ruled for the Warhol Foundation, but the Second Circuit Court of Appeals reversed, ruling that the four fair-use factors favored Ms. Goldsmith.¹⁹

The Supreme Court agreed to hear the case, but limited the issue on appeal: “Although the Court of Appeals analyzed each fair use factor, the only question before

this Court is whether the court below correctly held that the first factor, ‘the purpose and character of the use, including whether such use is of a commercial nature or is for nonprofit educational purposes,’ 107(1) weighs in Goldsmith’s favor.”²⁰

The Court explained that “the first fair use factor...focuses on whether an allegedly infringing use has a further purpose or different character, which is a matter of degree, and the degree of difference must be weighed against other considerations, like commercialism.”²¹ “The [Copyright] Act’s fair use provision...‘sets forth general principles, the application of which requires judicial balancing, depending on relevant circumstances.’”²² The first fair-use factor, according to the Court, addressed the issue of substitution, “‘copyright’s bete noire.’ The use of an original work to achieve a purpose that is the same as, or highly similar to, that of the original work is more likely to substitute for, or ‘supplant’ the work.”²³ The first factor, which is just one of the statutory factors to be considered in fair use, “asks ‘whether and to what extent’ the use at issue has a purpose or character different from the original....The larger the difference, the more likely the first factor weighs in favor of fair use. The smaller the difference, the less likely.”²⁴

The Court then turned to an examination of the term “transformative use,” which is a use with a further purpose or different character. “[T]ransformativeness is a matter of degree.” The Court observed that the word transform does not appear in 107, but does appear in defining derivative works. “The statute defines derivative works, which the copyright owner has ‘the exclusive righ[t]’ to prepare, 106(2), to include ‘any other form in which a work may be recast, transformed, or adapted,’ 101. In other words, the owner has a right to derivative transformations of her work....To be sure, this right is ‘[s]ubject to’ fair use....**But an overbroad concept of transformative use, one that includes any further purpose, or any different character, would narrow the copyright owner’s exclusive right to create derivative works.**”²⁵ The Court also noted that “[a] use that shares the purpose of a copyrighted work...is more likely to provide ‘the public with a substantial substitute for matter protected by the [copyright owner’s] interests in the original wor[k] or derivatives of [it].”²⁶ Also weighed against whether the use has a further purpose or different character is “the fact that a use is commercial as opposed to nonprofit....”²⁷

Summarizing how it will apply the first fair-use factor, the Court stated that “[i]f an original work and a secondary use share the same or highly similar purposes, and the secondary use is of a commercial nature, the first factor is likely to weigh against fair use, absent some other justification.”²⁸

The Court limited its analysis to Warhol’s licensing of Orange Prince to Conde Nast. At the outset, the Court observed that both Ms. Goldsmith’s and Warhol’s use of the photograph was the same, that is, the use of portraits of Prince used in magazines to illustrate stories about Prince.²⁹ Accordingly, Warhol’s licensing of the Orange Prince shared the same objectives of Goldsmith’s photograph, “even if the two were not perfect

substitutes.”³⁰ And both uses were commercial. Taken together, these two elements, the Court opined, weighed against a finding of fair use.³¹

“Transformativeness,” the Court observed, could in some instances outweigh its commercial character, but not here.³² Relying on its own precedent, the Court observed that the first fair-use factor “cannot be read to mean that 107(1) weighs in favor of any use that adds some new expression, meaning, or message.”³³ “Otherwise, ‘transformative use’ would swallow the copyright owner’s exclusive right to prepare derivative works.”³⁴ Meaning or message is “simply relevant to whether the new use served a purpose distinct from the original, or instead superseded its objects. That was, and is, the ‘central’ question under the first factor.”³⁵ The Court identified the standard for applying the first fair use factor: “Whether the purpose and character of a use weighs in favor of fair use is...an objective inquiry into what use was made, *i.e.*, what the user does with the original work.”³⁶

Applying this objective standard, the Court ruled that Warhol’s commercial use of the portrait of the Orange Prince to illustrate a magazine about Prince, even though that portrait portrays Prince “somewhat differently from Goldsmith’s photograph (yet has no critical bearing on her photograph), that degree of difference is not enough for the first factor to favor [Warhol], given the specific context of the use.”³⁷ The Court also noted that “[Warhol] offers no independent justification, let alone a compelling one, for copying the photograph, other than to convey a new meaning or message. As explained, that alone is not enough for the first factor to favor fair use.”³⁸ In further support, the Court observed that copying can help to convey a new meaning or message. But that is not enough to satisfy the first fair-use factor. “Nor does it distinguish [Warhol] from a long list of would-be-fair users: a musician who finds it helpful to sample another artist’s song to make his own, a playwright who finds it helpful to adapt a novel, or a filmmaker who would prefer to create a sequel or spinoff, to name just a few.”³⁹

Finally, the Court cautioned that the four fair-use statutory factors may not be isolated from one another. “All are to be explored, and the results weighed together, in light of the purposes of copyright.”⁴⁰ Even though the Court focused only on the first factor, it affirmed the judgment of the Court of Appeals that Warhol had failed to prove fair use.⁴¹

Justice Gorsuch wrote a concurring opinion which was based, in part, on an interpretation of two provisions in the Copyright Act. On the one hand, a copyright holder has the right to create derivative works that “transform” or “adapt” the original work.⁴² Therefore, claiming that a later user “transformed” the work by endowing it with a new message or aesthetic cannot automatically mean that the subsequent use is fair. “To hold otherwise would risk making a nonsense of the statutory scheme....”⁴³

Justice Kagan wrote a harsh dissent. The Justice wrote that the majority had dramatically altered the law of fair use. No longer is the issue “[d]oes the work add something new, with a further purpose or different character, altering the [original] with

new meaning or expression.”⁴⁴ Instead, “[a]ll that matters is that Warhol and the publisher entered into a licensing transaction, similar to one Goldsmith might have done. Because the artist has such a commercial purpose, all the creativity in the world could not save him.”⁴⁵

In sum, this U.S. Supreme Court ruling restricts the scope of fair use, because a transformative use is no longer enough.

Citing *Warhol*, Courts Have Noted A Material Change In The Fair Use Standard

Applying the holding in *Warhol*, several courts have underscored the importance of comparing the specific uses of the original work and the copy in applying fair use. Discussing *Warhol*, the U.S. District Court for the Central District of California in *Sedlik v. Drachenberg*, 2023 U.S. Dist. LEXIS 183184 (C.D. Cal. 2023) noted that “[t]he Court agrees with the parties that there has been a material change in [the] controlling law [of fair use]....[that] changes the Court’s analysis.” *Id.* at *4. While no case has directly addressed the new fair use standard in the context of Generative AI, those cases have cited *Warhol* for the proposition that “new meaning or message is not sufficient....Instead meaning or message [is] simply relevant to whether the new use served a purpose distinct from the original, or instead superseded its objects. That was, and is the central question under the first factor [of a fair use analysis.]” *Id.* at *6-7, quoting *Warhol*, 143 S. Ct. at 1284 (2023). “A court should not attempt to evaluate the artistic significance of a particular work[.]” 2023 U.S. Dist. Lexis at *7, quoting *Warhol*, 143 S. Ct. at 1283; *see also*, *Cramer v. Netflix*, 2023 U.S. Dist. LEXIS 165510 (W.D. Pa. 2023) (“The Supreme Court instructs that an artist’s stated or perceived intent does not dictate whether a work is transformative, but rather ‘[w]hether the purpose and character of a use weighs in favor of fair use is, instead an objective inquiry into what use was made, *i.e.*, what the user does with the original work.’”) *Id.* at *19-20, quoting, *Warhol* at 1284; *Larson v. Perry*, 2023 U.S. Dist. LEXIS 163833 (D. Mass. 2023) (“To distinguish itself from a derivative use, then, a transformative use must serve a manifestly different purpose from the [work] itself.”) *Id.* at *31-32, quoting *Warhol* at 1274.

Generative AI, which, uses neural networks to rearrange excerpts of copyrighted literature and the visual arts to produce works for commercial purposes will have difficulty meeting the *Warhol* standard for fair use. The challenge for creators is to show that their works were actually used by LLMs, since companies such as OpenAI no longer allow access to the information used to train LLMs. That is why tagging, discussed later in this paper, is so important.

Some who claim training LLMs is fair use rely, in large part, on two federal cases: *Google LLC v. Oracle Am., Inc.*, 141 S. Ct. 1183 (2023) and *Authors Guild v. Google, Inc.*, 804 F. 3d 1183 (2d. Cir. 2014). Neither case supports a fair use defense regarding training LLMs. In *Google LLC v. Oracle Am., Inc.*, the U.S. Supreme Court ruled that Google’s unauthorized copying of a portion of a computer code called Java SE was fair use. 141 S.

Ct. at 1197. The code was a user interface, allowing programmers to control task-performing computer programs “via a series of menu commands.” *Id.* at 1201. The Court observed that “[n]o one claims the decisions about what counts as a task are themselves copyrightable.” *Id.* In *Google*, the code did not make use of copyrighted works. Training LLMs, by contrast, does make use of reams of copyrighted work.

In *Authors Guild v. Google, Inc.*, the plaintiffs, authors of published books protected by copyright, sued Google for making digital copies of tens of millions of books that were submitted to it by major libraries. Google scanned the copies and established a publicly available search function, which allowed an internet user to search without charge whether the book contained a specified word or term, and also observe “snippets” of text containing the searched-for terms. The court ruled that making a digital copy to provide a search and snippet function was a “transformative use.” 804 F. 3d at 207. Discussing the first fair use factor, the court noted that Google’s making a copy of plaintiff’s books for enabling a search for identification of books containing a term of interest was transformative. 804 F.3d at 216. This search function served a different purpose from that of the plaintiff’s books, which were sold commercially to be read in full, either in an actual paper-bound book or on the internet. Accordingly, it meets the *Warhol* test of the first fair use standard, a different purpose. In addition, Google’s search function also satisfied the fourth fair use factor which focuses on whether the copy brings to the market a competing substitute for the original, or its derivative. The Court ruled that the search and snippet function “does not give searchers access to effectively competing substitutes.” 804 F.3d at 224.

Unlike the Google search and snippet functions in *Author’s Guild*, LLMs produce works based on copyrighted material for commercial purposes, the very same purpose as the copyrighted works. LLMs do not produce search terms or snippets; they combine copyrighted works into a final product that bears no resemblance to search terms or snippets. Training LLMs also fails the fourth factor, because they produce competing substitutes, or derivatives. Accordingly, under *Warhol* and factor four, fair use would not apply.

**The Copyright Office Must Issue a New Policy Regarding
AI-Produced Works and Fair Use In Light of the U.S. Supreme Court’s
Decision In *Andy Warhol Found. for the Visual Arts, Inc. v. Goldsmith*.**

The Copyright Office has a Fair Use Index, a project undertaken by the Office of the Register in conjunction with the Intellectual Property Enforcement Coordinator. “The goal of the Index is to make the principles and application of fair use more accessible and understandable to the public by presenting a searchable database of court opinions....”⁴⁶ The index tracks judicial decisions “to help lawyers better understand the types of uses

courts have previously determined to be fair – or not fair.”⁴⁷ The Copyright Office does not offer opinions on fair use to members of the public.⁴⁸

The Index summarizes the four statutory factors that courts take into account to determine fair use. With respect to the first factor, the Office states that “[c]ourts look at how the party claiming fair use is using the copyrighted work, and are more likely to find that nonprofit educational and noncommercial uses are fair.”⁴⁹ This does not mean that “all commercial uses are not fair; instead, courts will balance the purpose and character of the use” against the other three statutory factors.⁵⁰ “Additionally, ‘transformative’ uses are more likely to be considered fair. Transformative uses are those that add something new, with a further purpose or different character, and do not substitute for the original use of the work.”⁵¹

The Copyright Office concluded its discussion of fair use by cautioning that “other factors may also be considered by a court in weighing a fair use question, depending on the circumstances. Courts evaluate fair-use claims on a case-by-case basis, and the outcome of any given case depends on a fact-specific inquiry. This means that there is no formula to ensure that a predetermined percentage or amount of work – or specific number of words, lines, pages, copies – may be used without permission.”⁵²

The Copyright Office updates the Index periodically, but at this point makes no mention of AI. Nor does the Index include a reference to the Supreme Court’s decision in *Warhol*. Doubtless the Copyright Office will amend its Index to provide new guidance on both, but how will it take into account the competing goals of stakeholders on the use of copyrighted works to train LLMs, and the output from LLMs? And what, if any, impact will *Warhol* have on the Copyright Office’s guidance on copyright infringement and fair use in this context? This article will offer a few observations.

The Copyright Office is unlikely to take a position on whether the action of an LLM in scraping data from the internet is copyright infringement. To be sure, LLMs scrape copyrighted works, but infringement also requires “substantial similarity” between the copyrighted material and the copy.⁵³ In any case, because infringement is decided on a case-by-case basis, the Copyright Office is not in a position to offer guidance on this issue, beyond a simple statement of existing law.

But whether AI companies can claim fair use for outputs produced by LLMs that have scraped the internet is a different matter, largely because of the Supreme Court’s decision in *Warhol*. The Copyright Office should address the applicability of fair use in this context, relying on the Court’s elevation of commercial over transformative use. The Court’s words on this issue are clear:

The dissent’s conclusion—that whenever a use adds new meaning or message, or constitutes creative progress in the opinion of critic or judge, the first fair use factor weighs in its favor—does not follow from its basic premise. Fair use instead strikes a balance between original works and secondary uses based in part on objective indicia of the use’s purpose and

character, including whether the use is commercial and, **importantly, the reasons for copying.**⁵⁴

The Court noted that copyright protection includes the right to prepare derivative works that transform the original. To ensure the protection of those rights, the Court went on to say that “[t]he use of a copyrighted work may nevertheless be fair if, among other things, the use has a purpose and character that is sufficiently distinct from the original. **In this case, however, Goldsmith’s original photograph of Prince, and [Warhol’s] copying use of that photograph in an image licensed to a special edition magazine devoted to Prince, share substantially the same purpose, and the use is of a commercial nature.**”⁵⁵

The Copyright Office must now issue new fair use guidelines underscoring that if both the copyrighted work and the copy have substantially the same purpose and the use is commercial, the defense of fair use is unlikely to apply. Even though the *Warhol* decision did not address AI, its ruling will nevertheless have a major impact on AI outputs, which are largely commercial. For example, if a novelist could show that a writer and an AI company infringed her work by using a similar plotline in a new novel, with a few superficial changes such as time and place (much like coloring a silkscreen of the Prince photograph), the AI company could not assert the defense of fair use, because both works were created for substantially the same use and for commercial purposes – selling a mystery novel about lawyers.

To be sure, proving that a writer’s work was the basis of an AI-created work would be a difficult task, given the enormous store of material in an LLM. One creative solution offered in the Listening Sessions was transparency, which could be accomplished by “tagging” works scraped by LLMs to indicate authorship, and then listing the tagged items used in the AI output. This would permit the AI output to be compared to the input, allowing a determination of infringement or fair use. Tags could give creators leverage to claim a license, something they do not now have. More important, it introduces fairness into the process by allowing both creators and AI companies to have the information they need to strike a fair bargain. Though this may increase the cost of the technology, given the enormous sums invested in this technology already with far more to come,⁵⁶ it only seems fair to give creators an opportunity to share in some of the revenue that made these AI-generated works possible.

Tags would also be consistent with the *Warhol* Court’s concern with protecting the derivative rights of creators: “[Copyright] protection includes the right to prepare derivative works that transform the original.”⁵⁷ Once creators understood how their work was used to generate AI output, they would be in a better position to determine if their derivative rights had been violated.

How would tagging work in practice? Let’s add a few details to the example of the writer who used AI to create his work. Assume the narrative is about a young lawyer who discovers fraud and other criminal behavior in the law firm that just hired him, and

encounters threats of physical harm and death when attempting to expose that criminality. But the law firm is now in New York and the young lawyer is an experienced trial lawyer, fresh out of the Office of the U. S. Attorney for the Southern District of New York, who wants to make his mark in private practice. Assume the writer gave this plotline to ChatGpt. Tagging would create a list of copyrighted works written by such acclaimed novelists as John Grisham, which were scraped by the LLM for use by the writer. The Copyright Office could then require that the list be submitted for filing in its new Recordation System, which, in turn, would notify Mr. Grisham and other novelists that they were on the list.⁵⁸ They could then determine how, if at all, their works were copied by the writer. This would give novelists such as Mr. Grisham an effective tool to protect their rights, and to hold AI companies accountable.⁵⁹ And under *Warhol*, fair use is likely not available as a defense.

Listing works used to produce an AI product would also benefit AI companies accused of infringement, because they would be in a better position to analyze the validity of an infringement claim. In short, transparency would benefit both parties to an infringement claim. This could result in fewer lawsuits, and more settlements.

Conclusion

The Copyright Office does not litigate cases, but its guidelines and regulations, especially in a new area such as AI, can implement the constitutional mandate to protect the rights of creators and inventors and inspire further contributions.⁶⁰ It is well within the Copyright Office's statutory authority to protect the rights of authors by attaching a tag to protect their copyrighted works. Tagging would also encourage new creative works by sending a message that these works will be protected. Placing tags on data clawed from the internet by LLMs would be an acceptable compromise to meet the interests of the various stakeholders. While the Copyright Office has no authority to make tags retroactive, courts that are currently litigating infringement and fair use issues do have the power to impose this remedy as part of an overall solution to protect the rights of authors.

President Biden's Executive Order begins by stating that "Artificial intelligence (AI) holds extraordinary potential for both promise and peril."⁶¹ His sweeping order is intended to keep a close eye on AI companies to prevent a variety of harms including misinformation by requiring "watermarks" on AI content. Watermarking and tagging serve the same purpose: transparency. Simply put, tagging would benefit authors, AI companies and the general public. To be sure, Congress, and most likely the Supreme Court, will have the final say, but the Copyright Office is in a unique position to offer important guidance on the use of AI as it becomes more central to business and culture.

¹ In this article, the author will refer to generative artificial intelligence as “AI” for ease of reference.

² *Andy Warhol Found. For the Visual Arts, Inc. v. Goldsmith*, 143 S. Ct. 1258, 1273-74 (2023), citing 17 U.S.C. § 701(1).

³ Ross Anderson, *Inside the Revolution at Open AI*, *The Atlantic* 52, 55-56 (Sept. 2023).

⁴ October 30, 2023, Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence. Para. 4.6(b). The Order is designed to protect against the potential risks of AI systems by, among other things, requiring AI developers to share safety test results with the U.S. Government; establishing standards and best practices for detecting AI-generated content and authenticating official content; advancing equity and civil rights; developing principles to mitigate the harms and maximize the benefits of AI for workers; and promoting innovation and competition in the AI space. See, October 30, 2023 FACT SHEET: President Biden Issues Executive Order on Safe, Secure, and Trustworthy Artificial Intelligence. Thus, President Biden aims to chart a middle path, allowing AI development to continue largely undisturbed while keeping a close eye on the AI industry. See, Kevin Roose, *Administration Aims for Balance with A.I. Order*, *NY Times* at B6 (Nov. 1, 2023).

⁵ Vol. 88, No. 51, *Federal Register* at 16190 (March 16, 2023).

⁶ *Id.*

⁷ *Id.*

⁸ Copyright Office Transcript of Proceedings In the Matter of Copyright and Artificial Intelligence Literary Works, including Software, Listening Session at 50 (April 19, 2023).

⁹ *Id.* at 12-14.

¹⁰ *Id.* at 23, 25.

¹¹ *Id.* at 29-30.

¹² *Id.* at 43-45.

¹³ *Id.* at 78.

¹⁴ *Id.* at 91-92.

¹⁵ *Id.* at 130-131.

¹⁶ *Andy Warhol Found. for the Visual Arts, Inc. v. Goldsmith*, 143 S. Ct. at 1266-71.

¹⁷ *Id.* at 1269.

¹⁸ *Id.* at 1270.

¹⁹ *Id.* at 1271-72.

²⁰ *Id.* at 1272-73. The other fair use factors are: “the nature of the copyrighted work; the amount and substantiality of the portion used in relation to the copyrighted work as a whole; and the effect of the use upon potential market for or value of the copyrighted work.” 17 U.S.C. §§ (2), (3), (4).

²¹ *Id.* at 1273, citing *Campbell v. Acuff-Rose Music, Inc.*, 510 U.S. 569, 579 (1994).

²² *Id.* at 1274, citing *Google LLC v. Oracle America, Inc.*, 593 U.S.____, 141 S. Ct. 1183, 209 L. Ed. 2d 311 (2021) (slip op. at 15).

²³ 143 S. Ct. at 1274.

²⁴ *Id.* at 1275, citing *Campbell*, 510 U.S. at 579.

²⁵ *Id.* at 1275 (emphasis added).

²⁶ *Id.* at 1276.

²⁷ *Id.* at 1276.

²⁸ *Id.* at 1277.

²⁹ *Id.* at 1278.

³⁰ *Id.* at 1279.

³¹ *Id.* at 1280.

³² *Id.* at 1280.

³³ *Id.* at 1282, citing *Campbell*.

³⁴ *Id.*

³⁵ *Id.* at 1283-84.

³⁶ *Id.* at 1284.

³⁷ *Id.* at 1284-85.

³⁸ *Id.* at 1285-86.

³⁹ *Id.* at 1286.

⁴⁰ *Id.* at 1287.

⁴¹ *Id.*

⁴² *Id.* at 1289, citing 17 U.S.C. §§ 101, 106(2).

⁴³ *Id.* at 1289.

⁴⁴ *Id.* at 1292, citing *Campbell v. Acuff-Rose Music, Inc.*, 510 U.S. 569, 579 (1994); see *Google*, 593 U.S. at __ (slip op. at 24).

⁴⁵ *Id.*

⁴⁶ U.S. Copyright Office Fair Use Index at 1.

⁴⁷ *Id.*

⁴⁸ *Id.*

⁴⁹ *Id.* at 2.

⁵⁰ *Id.*

⁵¹ *Id.*

⁵² *Id.*

⁵³ *Craft & Simth, LLC v. E.C. Design, LLC*, 969 F.3d 1092, 1101 (10th Cir. 2020).

⁵⁴ 143 S. Ct. at 1287 (emphasis added).

⁵⁵ *Id.* (emphasis added).

⁵⁶ Deepa Seetharaman & Berber Jin, *OpenAI Seeks New Valuation of Up to \$90 Billion in Sales of Existing Shares*, Wall Street Journal, Sept. 26, 2023.

⁵⁷ *Warhol*, 143 S. Ct. at 12.

⁵⁸ The Recordation System is designed to accept any “document pertaining to a copyright.” 17 U.S.C. § 512 (a)(2). Such a document could include the digital address of an author whose copyrighted work appears on the list of scraped documents. The Recordation System has a centralized messaging center, which would automatically

notify the listed authors. See Recordation System at copyright.gov. Such a notification system is within the statutory guidelines of the Recordation System, and could be implemented using AI software.

⁵⁹ Tagging would provide creators with significant proof of infringement and derivative use. This might very well lessen or even eliminate the cost of litigation, and lead to faster settlements. In addition, identifying a large number of copyright infringements could cause Congress to enact protections, such as authorizing collective bargaining and licensing. Recently, OpenAI and Microsoft have offered to cover the legal costs of the business users of ChatGpt and Microsoft Copilot, respectively, accused of copyright infringement. It remains to be seen whether these offers will increase or decrease infringement litigation.

⁶⁰ “The Copyright Office has statutory authority to issue regulations necessary to administer the Copyright Act.... The Copyright Office has authority to interpret the Copyright Act, and its interpretations of the Act are due deference.” *Motion Picture Ass’n of Am. v. Oman*, 750 F. Supp. 3, *6 (D.D.C. 1990) (citations omitted).

⁶¹ Executive Order of October 30, 2023 at Section 1. The private sector has also weighed in calling for transparency. The Data & Trust Alliance, a nonprofit group made up of two dozen large companies, including American Express, Humana, IBM, Pfizer and others, has developed standards for describing “the origin, history and legal rights to data. The standards are essentially a labeling system for where, when and how data was collected and generated, as well as its intended use and restrictions.” Steven Lohr, *Companies Find Way to Spot Trustworthy A.I. Data*, New York Times, December 1, 2023.