Before the

# U.S. COPYRIGHT OFFICE, LIBRARY OF CONGRESS

**Artificial Intelligence and Copyright**
Docket No. 2023-6

**Comments of Electronic Frontier Foundation**
October 30, 2023

Submitted by:
Kit Walsh
Director of AI and Access-to-Knowledge Legal Projects
Electronic Frontier Foundation
815 Eddy Street
San Francisco, CA  94109
Telephone: (415) 436-4333
kit@eff.org

The Electronic Frontier Foundation (EFF) is the leading nonprofit organization defending civil liberties in the digital world. For more than 30 years, EFF has represented the public interest in ensuring that law and technology support human rights and innovation. In the United States and abroad, we work to ensure that copyright policy, legislation, and practice appropriately serve the public interest by striking a just balance between interests of professional creators and innovators and the general public.

As requested, the form of EFF's comments will be specific responses to the Copyright Office's August 30, 2023, Notice of Inquiry.

**1. As described above, generative AI systems have the ability to produce material that would be copyrightable if it were created by a human author. What are your views on the potential benefits and risks of this technology? How is the use of this technology currently affecting or likely to affect creators, copyright owners, technology developers, researchers, and the public?**

Generative AI has proven a valuable technology for producing low-quality text such as machine translations and as an assistive tool for creators working in visual and auditory media. For example, the ability to translate a web page or email instantly and for free has opened up a range of learning and communication that would never occur if the alternative would be to hire a human being. At the same time, higher-quality translations require a human level of understanding to capture context, nuance, and to localize expression that would be unintelligible if translated literally.

While the development of technology is unpredictable, there are reasons to believe that current generative AI technologies are limited by their very nature to a relatively low quality of output. Emulating common patterns in a large corpus of works is a recipe for rehashing the common themes of a genre or producing the most archetypal rendering of an object, but not for breaking new artistic ground except by the surprising random collisions unique to machine learning – or in the hands of a creative person.

We believe it is possible that this automation will displace some demand for the low end of some creative labor markets, in a way that is comparable to other forms of labor being replaced by automation. Both creators and other workers are reasonably concerned about losing the ability to make a living with their labor, an important social and policy challenge that must be addressed with broader and more flexible tools than copyright law. *See generally* Corynne McSherry, "Generative AI Policy Must Be Precise, Careful, and Practical: How to Cut Through the Hype and Spot Potential Risks in New Legislation," (July 7, 2023), https://www.eff.org/deeplinks/2023/07/generative-ai-policy-must-be-precise-careful-and-practical-how-cut-through-hype.

**4. Are there any statutory or regulatory approaches that have been adopted or are under consideration in other countries that relate to copyright and AI that should be considered or avoided in the United States? How important a factor is international consistency in this area across borders?**

The United States enjoys the robust protection for free expression embodied in the First Amendment, and the limits on the scope of copyright enforcement that protect First Amendment values, such as fair use and the idea/expression dichotomy. Approaches in other jurisdictions should not be copied here without taking into account differences in the underlying legal protections for speech.

**5. Is new legislation warranted to address copyright or related issues with generative AI? If so, what should it entail? Specific proposals and legislative text are not necessary, but the Office welcomes any proposals or text for review.**

Existing copyright law is sufficiently flexible to evaluate the copyrightability and infringement questions arising from generative AI. In fact, multiple district courts around the country are evaluating these questions in various cases and will soon provide additional guidance on its application. As for other issues, such as impact on creative labor markets, those social and policy challenges must be addressed with broader and more flexible tools than copyright law.

**8. Under what circumstances would the unauthorized use of copyrighted works to train AI models constitute fair use? Please discuss any case law you believe relevant to this question.**

Training an AI model may implicate the reproduction right when works are copied (e.g., from the internet) to be stored together as a training data set. Training then consists of analyzing those works in order to store information about the patterns that arise when they are evaluated in concert with the other elements of the data set. These observations are then stored as a model that can be used to generate new works.

Like other forms of intermediate copying, studying copyrighted works to form observations about patterns arising from a large corpus is likely to be a fair use in most circumstances. The freedom to analyze copyrighted works is part of the traditional contours of copyright, specifically both the idea/expression dichotomy and Section 107's favored purposes: commentary, criticism, and educational use. With respect to the former, it is crucial to note that a machine learning process is not limited to, or even primarily, identifying patterns that reflect copyrightable elements of an original work. Rather, much of its "learning" is focused on non-copyrightable properties of the subject being depicted (limes being green, trees having a vertical trunk with branches, cats having four legs, etc.) or *scenes-a-faire* that are free for all to use (formulaic premises and plot twists in language models, archetypal uses of visual language and color to convey mood). Put another way, since current AI technologies

commonalities within the datasets, they are necessarily looking elements that are least likely to be original and copyrightable.

Multiple federal courts are considering the question of liability in various real-world contexts, rather than in the abstract. Congress should allow the courts to do their job and interpret the law as it stands before considering legislation that would alter the legal landscape of this still-emerging technology, or making assumptions about how it will develop.

**8.1. In light of the Supreme Court's recent decisions in *Google v. Oracle* America and *Andy Warhol Foundation v. Goldsmith,* how should the "purpose and character" of the use of copyrighted works to train an AI model be evaluated? What is the relevant use to be analyzed? Do different stages of training, such as pre-training and fine-tuning, raise different considerations under the first fair use factor?**

As discussed, a model is a collection of information about what has been expressed, specifically what patterns emerge when many works are compared. To the extent training the model requires reproducing copyrightable expression, that copying has a different purpose and message than the creative works being analyzed to generate these observations. And, as *Google v. Oracle* teaches, a use that facilitates the creation of new works is more likely to be fair. As in *Google*, a model can be used for a range of expression informed by user prompts, conveying messages devised by users.

**8.3. The use of copyrighted materials in a training dataset or to train generative AI models may be done for noncommercial or research purposes. How should the fair use analysis apply if AI models or datasets are later adapted for use of a commercial nature? Does it make a difference if funding for these noncommercial or research uses is provided by for-profit developers of AI systems?**

It has historically been common for innovations from noncommercial and research activities to be re-purposed in a commercial context. Noncommercial and research uses are favored in part because they tend to create public benefits and generate knowledge that is available to all, advancing the purposes of copyright law. When a use generates knowledge and resources available to all, it will rarely if ever be appropriate to penalize researchers and nonprofit users for the fact that others find their innovations valuable.

**8.4. What quantity of training materials do developers of generative AI models use for training? Does the volume of material used to train an AI model affect the fair use analysis? If so, how?**

For the range of real-world applications, the volume does not affect the analysis.

**8.5. Under the fourth factor of the fair use analysis, how should the effect on the potential market for or value of a copyrighted work used to train an AI model be measured? Should the inquiry be whether the outputs of the AI system incorporating the model compete with a particular copyrighted work, the body of works of the same author, or the market for that general class of works?**

Section 107 is clear that the inquiry concerns the effect on the market for "the copyrighted work" alleged to be infringed.

**9. Should copyright owners have to affirmatively consent (opt in) to the use of their works for training materials, or should they be provided with the means to object (opt out)?**

No. There is no opting out of fair use, as fair use does not require the permission of a copyright holder. Nor can one "opt-out" of the use of non-copyrightable ideas that may be reflected in copyrightable expression.

**9.1. Should consent of the copyright owner be required for all uses of copyrighted works to train AI models or only commercial uses?**

Consent of the copyright owner should only be required for otherwise infringing uses.

**9.3. What legal, technical, or practical obstacles are there to establishing or using such a process? Given the volume of works used in training, is it feasible to get consent in advance from copyright owners?**

It would not be feasible to seek authorization from every copyright owner, particularly since the elimination of formalities means that copyright attaches at fixation to all sorts of amateur creations not part of any market. The effect of requiring authorization would be to limit competition to companies that have their own trove of images or strike a deal with such a company, resulting in all the usual harms of limited competition (higher costs, worse service, security risks) as well as reducing the variety of expression used to train such tools *and* the expression allowed to users of such tools seeking to express themselves. *See generally* FTC, "Generative AI Raises Competition Concerns," (June 29, 2023), https://www.ftc.gov/policy/advocacy-research/tech-at-ftc/2023/06/generative-ai-raises-competition-concerns.

**12. Is it possible or feasible to identify the degree to which a particular work contributes to a particular output from a generative AI system? Please explain.**

In general, no, because the model has learned what patterns appear in the dataset and a single work does not constitute a pattern. A notable exception is when a single image is included in the training set many times. In that case, a model may treat the details of that image as a pattern, and in rare instances may even produce a near copy of the original ('memorization'). However, curators of training data try to avoid duplicates because memorization of specific documents is considered a bug, not a feature.

That said, even if it were possible to identify the general contribution of a particular work, it would be difficult to discern whether the particular elements of the model contributing to a given output were informed by copyrightable elements of that work, as opposed to non-original facts and ideas or scenes-a-faire.

**13. What would be the economic impacts of a licensing requirement on the development and adoption of generative AI systems?**

As discussed in response to Question 9.3, the effect of requiring authorization would likely be to limit competition to companies that already have their own trove of images or strike a deal with such a company, resulting in all the usual harms of limited competition (higher costs, worse services, less security) as well as reducing both the variety of expression used to train such tools *and* the expression allowed to users of such tools.

**14. Please describe any other factors you believe are relevant with respect to potential copyright liability for training AI models.**

This is a fast-moving and emerging area of technology. Recent advances in generating low-quality text and images have stoked both hype and anxiety about the possible future capabilities of generative AI, but such speculation is a poor basis for regulation. Congress should not act hastily, particularly in the realm of copyright, given that restrictions on use of copyrightable expression can last for life plus 70 years, and statutory damages are draconian.

**22. Can AI-generated outputs implicate the exclusive rights of preexisting copyrighted works, such as the right of reproduction or the derivative work right? If so, in what circumstances?**

In general, the outputs of AI models reflect a vanishingly small portion of the information contained in their training data. AI models learn patterns that appear in multiple elements of the training data, including scenes-a-faire or otherwise noncopyrightable elements. For example, photographs of cats in the training data help an AI image generator learn to associate 'cat' with fur, four legs, and a tail. The photographers do not hold copyright in that information.

In at least two circumstances, however, AI-generated outputs may implicate these exclusive rights. First, diffusion-based models that are trained with a dataset containing many copies of the same work can, in rare instances, 'memorize' the work and create something recognizably similar to it when prompted. This could implicate the right of reproduction or the derivative work right, though it would likely require the user to actively seek to cause that result. Makers of AI models seek to avoid duplication because it skews the model in the direction of the duplicated elements. Memorization is therefore not only rare, but something model makers have an incentive to avoid and something existing copyright law can handle.

Second, a user may have a work in mind that they desire to reproduce or from which they desire to create a derivative work, and may manage to prompt an AI to generate an infringing work. If so, any liability should reside with that user, given that they have done the illicit copying, and the infringed work may not even be in the training data of the AI tool; a person could coax a tool trained on entirely different inputs to make something substantially similar and thereby infringe.

**23. Is the substantial similarity test adequate to address claims of infringement based on outputs from a generative AI system, or is some other standard appropriate or necessary?**

Yes, and we do not recommend any change to the test. The substantial similarity test is the traditional boundary on exercise of the right to control derivative works, and represents an essential balance between the incentive to create and the protection of creative works that build upon, reference, learn from, or otherwise could be said to 'derive' from other copyrighted works.

**25. If AI-generated material is found to infringe a copyrighted work, who should be directly or secondarily liable—the developer of a generative AI model, the developer of the system incorporating that model, end users of the system, or other parties?**

It should be quite rare for the developer of a model or system to be secondarily liable, given that these technologies rarely produce infringing works and have substantial noninfringing uses. *See Sony Corp. of America v. Universal City Studios, Inc.*, 464 U.S. 417 (1984).

Since copying is an element of infringement, only a party who carried out the element of copying can be considered infringing. This standard rules out end users in situations where they were not using a tool to create a copy of a particular work, but happened to generate a potentially infringing work. It would be dangerous and unfair to create a category of unwitting infringers with no way to know that their conduct is unlawful and no reason to suspect it as long as generative AI technology continues to function as it does today.

**25.1. Do "open-source" AI models raise unique considerations with respect to infringement based on their outputs?**

By virtue of being understandable and customizable, truly open models are even more likely to have substantial non-infringing uses as they can be understood, used, and modified in more ways than a closed-source tool. Such models also promote desirable competition and innovation by reducing barriers to deploying and customizing AI technologies. Merely being marketed as "open," however, does not make it so. *See* Widder, David Gray and West, Sarah and Whittaker, Meredith, "Open (For Business): Big Tech, Concentrated Power, and the Political Economy of Open AI" (August 17, 2023), https://ssrn.com/abstract=4543807.

**28. Should the law require AI-generated material to be labeled or otherwise publicly identified as being generated by AI? If so, in what context should the requirement apply and how should it work?**

While labeling could have benefits, it largely speaks to considerations that have nothing to do with copyright or the purview of the Copyright Office, such as concerns specific to misinformation. Any such proposals must be narrowly tailored to the harms they seek to address, and in particular avoid several potential pitfalls. Any requirement must not place such a burden on those making or using generative AI that it forecloses nonprofits or small companies from participating in such innovation. Second, the labeling must not effectively create a new tracking mechanism at the expense of Americans' rights to anonymous expression and reading.

**31. Should Congress establish a new federal right, similar to state law rights of publicity, that would apply to AI-generated material? If so, should it preempt state laws or set a ceiling or floor for state law protections? What should be the contours of such a right?**

Only if that right is sharply limited and offers the added benefit of preempting overbroad state provisions.

State rights of publicity have far overstepped their legitimate bounds in many jurisdictions. While a privacy-based right of publicity to avoid false commercial endorsements is appropriate, some states have transformed this privacy interest into a quasi-property right that can be used to control or chill creative speech such as a robot parody of a celebrity, depictions of real-world events or alternate scenarios involving a person (such as sports games depicting the name and appearance of the players on the various teams), jokes involving the person's name (like "Here's Johnny!" on a portable toilet), a car associated with a person, and even depictions of a famous person. *See Motschenbacher v. R.J. Reynolds Tobacco Co.*, 498 F.2d 821 (9th Cir. 1974) (use of modified plaintiff's race car); *Carson v. Here's Johnny Portable Toilets*, 698 F.2d 831 (6th Cir. 1983) (use of famous phrase associated with talk show host: "Carson's identity may be [wrongfully] exploited even if his name, John W. Carson, or his picture is not used."); Jennifer Rothman, The Right of Publicity: Privacy Reimagined for a

Public World, at 5, Harvard University Press, (2018). The courts have failed to adopt an effective legal standard for safeguarding First Amendment interests against the nebulous right to control speech that evokes one's identity.

The mission of copyright law, to foster the creation and dissemination of original creative works, is therefore in direct opposition to a broad federal right of publicity. A limited federal right of publicity that preempts broader state rights would be preferable.

In addition, Congress should clarify that the right of publicity sounds in privacy and is not "intellectual property" for purposes of Section 230 of the Communications Decency Act. As we have learned, when platforms must fend off expensive lawsuits to protect user speech, they are likely to cave to censorious demands. At best, we would see calls to expand fundamentally flawed systems like Content ID that regularly flag lawful content as potentially illegal and chill new creativity that depends on major platforms to reach audiences. Katharine Trendacosta, "Unfiltered: How YouTube's Content ID Discourages Fair Use and Dictates What We See Online," (Dec. 10, 2020), https://www.eff.org/wp/unfiltered-how-youtubes-content-id-discourages-fair-use-and-dictates-what-we-see-online. Today, some new creators avoid using any music at all lest they trigger a flag, however improper, and music criticism is famously underrepresented on platforms using such tools because it is effectively impossible. If creators were chilled to this degree when discussing or referencing any public figure, even those long dead, it would be a tremendous loss.

**32. Are there or should there be protections against an AI system generating outputs that imitate the artistic style of a human creator (such as an AI system producing visual works "in the style of" a specific artist)? Who should be eligible for such protection? What form should it take?**

As discussed in response to Question 23, it would undermine free expression and the goals of copyright law to suppress more original expression than that which is substantially similar to a copyrighted work of which it is derivative. This rule already forbids many uses that are primarily original, such as sequels or many uses of a works' setting or even individual characters. A greater degree of restriction on the public's permissible range of speech, particularly one as elusive as a "style," would undermine the cultural advancement at the core of copyright's goals.