Hello,

My name is Dylan Garity. I am a writer and editor by profession, and I cofounded and served as vice president of an independent publishing company for many years, where I worked with dozens of incredible authors and artists. I am submitting this comment in support of writers, artists, musicians, and all other people whose life's work is being stolen for the profit of corporations and tech developers under the guise of progress.

The facts are simple: Generative AI companies have copied millions of copyrighted artworks, writings, and recordings without consent or permission, and have used those copies to create tools and artificially generated works that directly compete with the original creators. This is a clear, blatant violation of fair use, especially in light of the US Supreme Court's recent ruling in *Warhol v. Goldsmith*.

The arguments submitted both here and in courts by companies such as OpenAI, Microsoft, Anthropic, and Meta demonstrate a cold understanding of what they have done and continue to do. They are trying to convince regulators and judges that because they have essentially stolen *everything*, they cannot be held accountable for stealing *anything*. The opposite is true—they must answer for every last piece.

These companies protest that they have invested billions into this technology and so cannot "go back." Yet why are their billions, risked by their own callous and careless choice, worth more than the life's work of millions, worth more than the lives and livelihoods of us all? They could've invested all the money in the world, and it would not give them the right to break the law.

They are right, however, that the longer they continue along this path, the harder it will be to rein in. That is why swift, sweeping action is necessary. The US Copyright Office is in a unique position to help with that action, and with affirming and expanding individual IP and likeness rights so as to prevent further criminal misuse.

On a legal and ethical level, the only appropriate resolution to this matter is the forced algorithmic disgorgement of all models that used copyrighted materials without consent in their training data, compensating every individual whose work or likeness was used without consent, and a robust, stringent system of review and regulation going forward for the development and release of any new generative model, to ensure that the training process and ultimate product is transparent, legal, and ethical. This technology can be useful for humanity, but not if it is based on a foundation of massive theft.

On the matter of if the outputs produced by these algorithms should be copyrightable, the answer is a resounding no. I applaud the US Copyright Office for its rulings to that effect to this point, and urge you to continue that strong stance. Any other choice would end the utility of copyright as we know it, and allow corporations and malicious individuals to completely destroy the creative professions and stifle future human creativity.

I have attached further responses to a number of the Copyright Office's specific questions as an addendum to this letter. Thank you for taking the time to read this, and for your important ongoing work and integrity in this matter.

Dylan

**Additional Responses to Select Questions:**

*2. Does the increasing use or distribution of AI-generated material raise any unique issues for your sector or industry as compared to other copyright stakeholders?*

Yes—I'd like to specifically emphasize the case of unethical third-party companies such as Sudowrite and Jasper, which use the underlying technology of major foundation models like GPT, plus extensive fine-tuning on further copyrighted materials, to offer a "writing partner" tool (which attempts to serve various functions traditionally fulfilled by both writers and editors).

In reality, what they are offering is a plagiarism machine, and a significant violation of my rights and the rights of millions of others. In Jasper, for example, you can specifically choose individual authors to base an output text upon, and it uses information copied and stolen from those authors' works in an attempt to do so. Sudowrite promises the ability to "write" a novel in a single weekend based on a fine-tuning system that involved prospective authors providing their books to the company under false pretenses. While neither product produces particularly impressive or consistent results at this moment, I have no doubt that they may continue to improve—especially if these companies are allowed to continue stealing copyrighted works for use in training and fine-tuning.

In both the creative writing and editing spaces, I expect many more specialized tools to continue cropping up, directly targeted at replacing humans and creating unfair competition, built on the backs of our pre-existing and protected work.

I want to be clear: I do not believe that "writing partner" tools based on generative AI should not be allowed to exist, period. While I have grave concerns for the future of humanity if generative AI is embraced at large as a viable replacement for human communication, creation, and connection, and that is an important larger ongoing discussion on a societal and governmental level, I can still see the potential utility of such tools with responsible development and usage. On a legal level, the concern is narrower and clearer: a "writing partner" tool that in any part of its training and development process involved the nonconsensual use of copyrighted materials should be disbanded and its developer held in violation of the law.

*3. Please identify any papers or studies that you believe are relevant to this Notice. These may address, for example, the economic effects of generative AI on the creative industries or how different licensing regimes do or could operate to remunerate copyright owners and/or creators for the use of their works in training AI models. The Office requests that commenters provide a hyperlink to the identified papers.*

I'd like to highlight two studies from this year, both of which showcase that contrary to the claims of Generative AI developers and tech companies, models do fundamentally retain (and illegally copy and output) significant portions of their copyrighted training data.

*Speak, Memory: An Archaeology of Books Known to ChatGPT/GPT-4:*
(https://arxiv.org/pdf/2305.00118.pdf)

*Scalable Extraction of Training Data from (Production) Language Models:*
(https://arxiv.org/abs/2311.17035)


These companies are trying to proclaim making an unauthorized copy of copyrighted material is negligible within the overall training process—but it is in many ways the entire point. Copyright laws are designed around human ability. I can read a book and be inspired by it without violating copyright, taking that experience in alongside the rest of my life. That is a part of the rich history of art and writing, and it is a process that comes with welcome human cognitive limitations. I must put more of myself into anything I then create (or else I am indeed liable for infringement if I merely attempt to copy). Even if these algorithms could just somehow "read" a text without that technical process involving theft/illegal copying, they do not have that same right, nor should they.

The idea that these "memorization" instances do not count as storing an illegal copy is absurd. For example, consider an algorithm with the instructed output of "always follow each word with the most likely next word. The most likely first word is 'Mr'. The most likely second word is 'and'. The most likely third word is 'Mrs'. The most likely fourth word is 'Dursley'"—and so on, for the entire text of *Harry Potter and the Sorcerer's Stone*.

That model would contain, of course, an illegal copy of the copyrighted book, despite the full text of the book never being directly present together—and these studies show that the current models are far more similar to that than their developers proclaim.


*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?*

International consistency should be secondary to human and ethical considerations. Though the US is unfortunately far behind and more piecemeal in terms of data regulation and privacy protection vs. for example Europe's GDPR, we are often looked to as a leader in global policy. Strict regulation of copyrighted data usage and against copyright of outputs is an opportunity for strong global leadership in that regard.

Another country creating policy that favors theft and big tech is not a justification for the US to do the same. "Oh other people will just steal things anyway," some argue—so does that mean we must strive to be the most unethical and even faster? We can instead shape not just US but global policy. Many other countries will follow in our strong, bold footsteps in this regard. US policy can and should be used to apply pressure on countries that may have gone down or be going down unethical and violating roads, to protect both the rights of Americans and those of other artists and authors worldwide.


*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.*

I am not a writer of legislative text so will leave the details to others in this case, but I think two major, intersecting pieces of legislation are necessary, in line with court rulings and other regulatory actions.

1) Legislation explicitly requiring consent for the use of *any* non-public-domain data for the training of generative AI models.

2) Legislation protecting one's likeness (appearance and voice) from copying, replication, and misuse. The horrific spread of unauthorized deepfakes and specifically deepfake pornography, including deepfake pornography of minors, shows that this cannot wait. Every individual should have a fundamental, legislated right to their own likeness and how it is used, particularly in the space of new technologies.


*7.3. Is it possible for an AI model to "unlearn" inferences it gained from training on a particular piece of training material? If so, is it economically feasible? In addition to retraining a model, are there other ways to "unlearn" inferences from training?*


I would like to highlight that I find the premise of the middle question here to be fundamentally misguided.

Given the ongoing development of generative tools that do *not* explicitly use stolen, copyrighted data (see Getty and Adobe), I do not believe it would be as impossible as companies say to retrain their models. I believe they are exaggerating to try to escape consequence. However, if they are not—if the fixing of this grave injustice would require so much work and money and time as to result in the dissolution and bankruptcy of companies like Meta and Google—then that is the price they and their owners and shareholders should have to pay. Millions of other individuals having their lives destroyed and livelihoods stolen should not be viewed as a necessary evil, while these companies are viewed as essentially immortal and unassailable.

It is possible to use legal, ethical data sources. While fundamental questions remain about if artists, photographers, and other creatives in the case of data used for something like Adobe Firefly intended their work to be able to be used in this fashion, and whether "use in any and all technologies now known or hereafter devised" should be legally upholdable contract language, that at least shows an intent to act legally and ethically, and that companies and developers using copyrighted works without permission think they can get away with blatant theft, and that the rights and futures of millions are worth less than cutting a few corners.


*7.4. Absent access to the underlying dataset, is it possible to identify whether an AI model was trained on a particular piece of training material?*

Yes, see the studies listed under Question 3. However, court- or legislatively-ordered access to the underlying dataset is likely necessary for finding if any *specific* piece of material was used.

I also have kept extensive personal documentation of various models' abilities to output, with minimal prompting and without permission, full poems and verbatim passages of my own writing, something they could not do without that writing being included in the training material.

*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?*

All three should be taken into consideration.


*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)?*

Affirmative consent is the only viable option. Opt-out does not come close to adequately addressing the issue, as it cannot be comprehensive, and seems primarily to be a false-compromise method for companies and developers to try to get away with ongoing theft.


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

Given the current abuse of training that was initially framed as or intended to be for research-only purposes but then co-opted or misused for commercial and profit-oriented purposes, consent should be required for all uses.

Limited exemptions for strictly regulated training related to medical advances and other "good of humanity" purposes could be a reasonable carveout worthy of healthy debate.


*9.4. If an objection is not honored, what remedies should be available? Are existing remedies for infringement appropriate or should there be a separate cause of action?*

Financial remedies are not enough. The forced disgorgement of any algorithm based on the nonconsensual use of copyrighted materials for training data is necessary, as well as additional punishments and restrictions for repeat violators (companies, teams, or individual developers).


*10.1. Is direct voluntary licensing feasible in some or all creative sectors?*

Again, I want to highlight and draw into question the very premise and supposition of "feasibility" here. Quite simply, if it is not feasible, then these works should not be used. The supposition that the training of these massive models is somehow necessary or a simple natural process that must go on one way or another, while human authorship and artistry are not, is on its face misguided and absurd. The question should be, "Is it feasible for authors and artists to be forced to give away their life's work for free and without consent to corporations that are intent on replacing them?" The answer to that is a clear no.

*15. In order to allow copyright owners to determine whether their works have been used, should developers of AI models be required to collect, retain, and disclose records regarding the materials used to train their models? Should creators of training datasets have a similar obligation?*

Yes, and yes.

*18. Under copyright law, are there circumstances when a human using a generative AI system should be considered the "author" of material produced by the system? If so, what factors are relevant to that determination? For example, is selecting what material an AI model is trained on and/or providing an iterative series of text commands or prompts sufficient to claim authorship of the resulting output?*

For the latter, there are no such circumstances without completely destroying the utility and purpose of copyright law. I may be the author of ten thousand iterative commands and prompts; that still does not make me the author of what the model spits out at the end of it all.

In terms of materials, that would only be potentially viable if all materials selected were already authored by the person in question. Even so, the ability of such a system to potentially generate millions of works at a rapid pace in a way that simply is not possible without this nonhuman technology would seem to serve only to encourage copyright trolls and a broad degradation of useful and usable copyright.

*20. Is legal protection for AI-generated material desirable as a policy matter? Is legal protection for AI-generated material necessary to encourage development of generative AI technologies and systems? Does existing copyright protection for computer code that operates a generative AI system provide sufficient incentives?*

That protection is utterly undesirable as a policy matter, serving only predatory corporations and other malicious actors at the expense of artists, authors, performers, and some of our dearest, most fundamental and universal human cultural values. That protection is not necessary to encourage development of those systems, as we can see they are actively being further developed despite the lack of that protection.

There is a necessary trade-off here, and it should not be immediately in favor of new and largely untested technology. People can still make art and write and perform without the use of generative AI tools, so the fundamental ability to create is not being taken away from anyone if those works are not copyrightable; however, if works produced using them are copyrightable, then the ability for people to do so in any public or financially viable way (including other people also using generative AI tools!) will be decimated, especially given the chilling legal implications.

The final question does raise a particular irony worth drawing attention to: If I were to take the copyrighted code of any of these generative AI companies and use it without permission in any remotely similar method to the way they are using my and other authors' copyrighted works, they would justifiably sue me to oblivion.

*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?*

It is not on its own adequate. The use on the input level must be taken into account regardless of the specific output. For example, developers and others can fine-tune models based on a specific author's body of work. If someone takes ten of my books and uses them to train a model with the specific goal of producing books like mine, and what it produces is not always exact passages of text from those books but texts similar enough in style and purpose as to be clearly competing with the originals, that must still be claimable as infringement. The process matters.

*24. How can copyright owners prove the element of copying (such as by demonstrating access to a copyrighted work) if the developer of the AI model does not maintain or make available records of what training material it used? Are existing civil discovery rules sufficient to address this situation?*

This is why upfront regulation and review before public release of any generative AI system is the best available option. Concerns over such regulation stifling innovation should not trump concerns over the ongoing theft of the life's work of millions.

*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?*

All of the above, but with much greater liability for the former, and significantly lesser liability for the latter. An end user who publishes a book after prompting a system to produce a book in the style of a particular author (using that author's stolen books), for example, should certainly be held liable, but the focus should be on the developers at all levels that are engaging in the initial and ongoing theft and thus making the end-user infringement possible.

*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?*

All AI-generated material that is shared publicly should be required to be labeled prominently. For images, a substantial watermark on the image should be required, and removal of that watermark should be illegal. For text, one option would be that above/before any artificially generated passage of text posted publicly, there should be a required disclaimer, and the intentional removal of that disclaimer should also be illegal. Various publicly available metadata options may also be relevant and useful for various forms of media to track the history and authenticity of a piece of media.

This kind of requirement may seem draconian, but holding together the social fabric necessitates unprecedentedly massive actions to meet an unprecedented moment. If no one can know what is real anymore and if there are no disincentives to mislead about what is real and what is not, that social fabric will fray beyond recognition.

*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?*

Yes. See my answers to question 5. To expand on that, a federal right of publicity/right of privacy in regard to AI training and AI-generated material should set a floor, not a ceiling, for state law protections.

*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?*

This is where the most stringent protections are necessary. Such style imitation is only possible, of course, with input/training data or fine-tuning data based on generally stolen copies of a creator's work. Comparisons to how a human artist is influenced by another's style are facetious—that is a celebrated part of the cultural history and tradition of human art, and is necessarily limited, while machine-generated style imitation is mechanical theft at scale.

No generative model, regardless of specific input, should be allowed to produce outputs "in the style of" any human creator without that creator's explicit, opt-in permission.

*34. Please identify any issues not mentioned above that the Copyright Office should consider in conducting this study.*

I would like to close with some final remarks, first drawing additional attention to two specific malicious ways in which copyrighted data becomes used for training without consent, both of which require an urgent combination of legislative, regulatory, and legal responses.

1) Social media platforms like Meta and Twitter/X have discussed scraping public posts for the creation of their models. They claim a nonexclusive but complete right to all content posted on their platforms. Given these platforms' centrality in society and commercial opportunity and their goal of serving as "public squares," forcing people to give away their rights in order to share their work with the public in such spaces is wrong. It also makes these platforms massive hubs for piracy, where people share millions of pieces of copyrighted artwork and writing without permission every day. While prior to the development of generative AI, this may have been a tacitly accepted state of things among copyright holders and publishers (better for people to share and credit and provide exposure to something than to waste significant time and money and energy trying to hunt it all down), that is no longer feasible given this competitive misuse. Stopping all such "forced consent" policies on major platforms in relation to generative AI, and holding those platforms accountable for their misuse of copyrighted materials, is a significant, immediate, legislative and regulatory need.

2) OpenAI and other companies claim the right to use and train models on all user inputs into systems such as ChatGPT. This, of course, means users can copy entire existing copyrighted texts or sets of copyrighted artwork into these tools, and the companies then claim the right to do whatever they want with them. Companies and developers should not be able to pass off the

responsibility for whether or not such material is stolen/copyrighted to the user. They must be held responsible for anything they ultimately use. If this means that they cannot train on user-submitted data period due to impracticality, that is their problem, not the problem of copyright holders. Users must also be held responsible and severely liable for intentional misuse of copyrighted materials, such as people fine-tuning image generators with hundreds of images from a particular artist without that artist's consent in order to spit out competing works. It is possible that when considering the scope of the potential damage here, the ability to "fine-tune" on the user end should not be legally allowable at all with generative AI models due to the rampant abuse of the process and the scale at which it allows infringement to occur. We have legal restrictions on the use of other tools that enable massive criminal enterprise, and that same kind of framework and thinking may be necessary to create policy here.

\*

Final remarks:

The longer government bodies and courts drag their feet on regulation and injunctions, the more developers and companies will try to claim that it's too late, that these models as they exist are already ingrained in too much. That will not make their argument any more valid. If my work is stolen and the product of that theft becomes so valuable that the people who stole it can't stand the thought of losing it, that is all the more reason for them to answer for it.

If these copyrighted materials are in fact not actually vital and important to developers and tech companies, then they should not use them. If they *are* so valuable and important that the companies will not willingly let them go, then that is a clear signal that they must gain consent and pay for that usage, not simply steal what they consider to be so valuable.

The ability of these tools to produce directly competitive works is nearly completely reliant on the original copyright violation. These tools cannot train on their own output—they collapse into nonsense in attempts to do so (https://arxiv.org/pdf/2307.01850.pdf). They fundamentally require human contribution, and that contribution cannot and must not be forced.

The fact that these models don't *just* function to output competing or infringing text but perform other tasks as well is also irrelevant. It is as if someone printed out a dozen pirated copies of my book and is claiming that because they can use the stack as a stepping stool, that negates the fact that they and all their friends can also read it for free and thus compete with my sales and my rights. Telling them they cannot do this is not telling them they can't have a stepping stool. It's saying they can't steal my property in order to have one rather than going and finding or making a perfectly nice one on their own.

These developers and companies could have worked to advance this technology without theft of copyrighted materials from the start. It may have been slower, and it may mean that a generative tool would not be able to do everything that it can do currently—but nor *should* it necessarily be able to. Whether the pursuit of better/stronger/faster machine learning models is a positive goal for humanity in general, however, is larger than the question at hand. Models can still be developed and function in many useful ways without theft and violation. But stopping these practices and creating a more robust system of protections for copyrighted work and for people's likenesses can prevent many of the most malicious uses.