



U.S. Copyright Office

**Notice of Inquiry on Artificial Intelligence & Copyright (Dkt. 2023–6)**

Reply Comments of Meta Platforms, Inc.

December 6, 2023

Meta Platforms, Inc. (“Meta”) appreciates this opportunity to submit reply comments on the United States Copyright Office’s (“Office”) Notice of Inquiry on Artificial Intelligence and Copyright. These reply comments focus primarily on a handful of points raised by rightsholders in their submissions to the Office.

**I. Introduction**

The comments submitted to the Office illustrate broad agreement on several critical issues. Commenters agree that Generative AI will have a transformative and beneficial impact on humanity’s creative potential and on fields like science and medicine. Commenters with expertise in AI development also agree on the mechanics of AI training relevant to the fair use analysis—including, for example, that the training process is essentially one of pattern recognition that does not exploit the expressive content of the material in the training corpus. (*See infra* Part II.)

Several commenters, however, insist that AI developers must rely entirely on licensed and public domain content for training purposes, citing a handful of Generative AI image and music models that were purportedly trained only on such content. But forcing developers to rely only on limited sets of licensed or public domain data will necessarily diminish AI systems’ ability to build a complete and accurate understanding of the concepts that underlie human communication, leading to less useful (and more biased) models and a weaker domestic AI industry. (*See infra* Part III.A.) Moreover, the licensed models to which these commenters point do not in any way establish that it is possible to create a state-of-the-art, general purpose AI model based entirely on data collected through licensing. (*See infra* Part III.B.)

Other commenters argue that content with greater commercial appeal provides more value to Generative AI models, and suggest that fair use does not permit AI models to learn from these “high quality” works. But “quality,” as used in the AI literature, refers not to commercial or aesthetic quality, but to the completeness of the data. In reality, a work with greater commercial appeal (*e.g.*, a well-written novel) is no more helpful to the development of an effective AI model as the same quantity of text from an internet comment board. The aesthetic “quality” of any piece of content is irrelevant, because AI training is an exercise in pattern recognition that in no way extracts or exploits the data’s expressive content. (*See infra* Part IV.)

Finally, commenters who oppose application of the fair use doctrine in this context universally overlook that doctrine’s essential function: to avoid “rigid application” of copyright laws that would undercut, rather than advance, the “Progress of Science and useful Arts.” Attempting to use copyright law to foreclose the development of a transformative and creativity-enhancing technology simply because it involves the creation of intermediate copies of copyrighted material—which never see the light of day, are often entirely temporary, and in no way interfere with rightsholders’ ability to sell or license their works in existing markets—is irreconcilable with this fundamental aspect of fair use doctrine. (*See infra* Part V.)

## II. There is Widespread Agreement About the Value of AI and the Technical Underpinnings of the Fair Use Defense

Numerous comments submitted to the Office underscored that artificial intelligence tools—like transformative tools of the past such as the printing press, the camera, or the computer—have the “potential to make everyone smarter and more capable.”<sup>1</sup> As TechNet noted in its comment, “Generative AI, in particular, has the potential to transform not only the creative industries, but other fields like software development, scientific research, healthcare, government administration, and education.”<sup>2</sup> Commenters noted how creators are using Generative AI tools to “create as many images as photographers created in the first 150 years of photography combined,”<sup>3</sup> translate their content to reach global audiences with ease,<sup>4</sup> and democratize visual effects in filmmaking.<sup>5</sup> And in the scientific industries, several comments highlighted how Generative AI is driving medical innovation,<sup>6</sup> transforming access to education,<sup>7</sup> mitigating climate change,<sup>8</sup> and improving everyone’s access to the world’s information.<sup>9</sup> As noted in Meta’s Initial Comment, Meta is a leader in making this revolutionary technology available to as many people as possible, as quickly as possible, through responsible open innovation.<sup>10</sup>

As multiple comments explained, Generative AI language models are trained “to identify relationships and patterns among words in a large dataset.”<sup>11</sup> Language models are trained to evaluate “the proximity, order, frequency, and other attributes of portions of words, called tokens,

---

<sup>1</sup> a16z Comment at 2.

<sup>2</sup> TechNet Comment at 1.

<sup>3</sup> OpenAI Comment at 2.

<sup>4</sup> *See id.*; *see also* Google Comment at 6 (discussing Google’s work building an Generative AI model “that will support the world’s 1,000 most spoken languages, bringing greater inclusion to billions of people historically marginalized or underserved communities all around the world”).

<sup>5</sup> OpenAI Comment at 3.

<sup>6</sup> *See* Google Comment at 5–6 (discussing how Google’s Generative AI tools like AlphaFold and Med-PaLM can “improve healthcare, including maternal care, cancer treatments, and tuberculosis screening”); *see also* a16z Comment at 2–3 (highlighting how “AI is driving medical innovation” in drug development, cancer diagnosis, and optimizing billing practices).

<sup>7</sup> *See, e.g.,* OpenAI Comment at 5 (discussing Duolingo’s use of GPT-4 is drastically increasing access to conversational practice in the language they are studying); a16z Comment at 3–4; *see also* Meta Initial Comment at 8 (discussing Straightlabs’s use of Meta’s AI technology to offer personalized mentoring through 3D avatars).

<sup>8</sup> *See* Google Comment at 6 (discussing how AI will play a key role in “mitigating and adapting to climate change: by tracking wildfire boundaries in real time; helping to reduce carbon emissions by decreasing stop-and-go traffic; and providing critical flood forecasts”).

<sup>9</sup> *See, e.g.,* Google Comment at 6 (discussing how Google’s Data Commons project “synthesizes publicly available data from government agencies and other authoritative sources into an open source, API-accessible knowledge graph available to everyone”).

<sup>10</sup> *See* Meta Initial Comment at 9.

<sup>11</sup> Google Comment at 5.

in its training data.”<sup>12</sup> The goal of these techniques is not to extract particular expressive content, but rather to extract syntactical, structural, linguistic, and other information from a corpus of works as a whole.<sup>13</sup> For that reason, numerous comments agreed that the quality of Generative AI models does not depend on the expressive elements of a particular piece of content, or the inclusion or exclusion of any individual piece of training data, but rather on the *quantity* and *diversity* of the training content as a whole.<sup>14</sup> Training data, in other words, is highly substitutable: as long as the model’s overall training corpus is large and diverse, the model will function just as effectively with or without any specific piece of content.

There was similarly widespread agreement that training AI models does not implicate the rights protected by copyright.<sup>15</sup> As some comments explained, that is because the “factual metadata and fundamental information that AI models learn from training data [is] not protected by copyright law” and that “[c]opyright law does not protect the facts, ideas, scènes à faire, artistic styles, or general concepts contained in copyrighted works.”<sup>16</sup> As TechNet correctly noted, AI models themselves cannot be fairly characterized as “derivative works” of the training data “because [AI] models do not ‘re-present [any] protected aspects of the original’ works to users.”<sup>17</sup>

Commenters also highlighted that Generative AI model training is squarely protected by the fair use doctrine.<sup>18</sup> Those commenters agree that the four fair use statutory factors, along with the many decades of case law interpreting them in the context of new technologies, confirm that it is not an infringement of copyright for a Generative AI model to “learn” by deriving statistical information from copyrighted texts, images, or other media.<sup>19</sup>

### **III. Forcing Developers to Train Only on Licensed or Public Domain Data Will Yield Weaker AI Models and a Weaker Domestic AI Industry**

Several commenters argued that AI developers should be required to train their models using only licensed or public domain data, and suggested that imposing such a requirement would not meaningfully impact the effectiveness or competitiveness of our domestic AI industry. Those commenters, for example, claim that “large-scale licensing of copyrighted works already happens in the marketplace today” and pointed to “a number of generative AI platforms that were trained

---

<sup>12</sup> *Id.* at 4.

<sup>13</sup> See OpenAI Comment at 6-7 (detailing the pre- and post-training process).

<sup>14</sup> See a16z Comment at 8.

<sup>15</sup> See, e.g., *id.* at 6 (“[G]enerative AI model training is a productive, non-exploitative use of training material . . . [that] does not exploit any protectable expression in any given work, and so it does not implicate any of the legitimate rightsholder interests that copyright law seeks to protect.”); see also TechNet Comment at 2–3.

<sup>16</sup> OpenAI Comment at 12.

<sup>17</sup> TechNet Comment at 3 n.6 (quoting *Authors Guild v. Google, Inc. (Google Books)*, 804 F.3d 202, 225–26 (2d Cir. 2015)).

<sup>18</sup> See, e.g., a16z Comment at 5–8; Google Comment at 8–11, OpenAI Comment at 12–14.

<sup>19</sup> See TechNet Comment at 5; see also a16z Comment at 5–8; Google Comment at 8–11, OpenAI Comment at 12–14.

on entirely licensed and/or public domain content.”<sup>20</sup> Some commenters, for example, pointed to two image-generation models—in particular, Adobe’s Firefly model and a model currently being developed by Nvidia and Getty Images—as examples of such Generative AI platforms.<sup>21</sup> These arguments are misguided, as discussed below.

#### A. Limiting the Scope of Training Data Will Diminish AI Development

First, these arguments reflect a misunderstanding of the nature of AI models. AI systems are only useful to the extent that they can develop a complete and accurate understanding of the concepts that underlie human communication. All communication—whether between friends and co-workers, or between a person and an AI system—relies on shared understanding of a range of concepts, from grammatical rules, to vocabulary, to cultural tropes, to common knowledge. An AI system that does not understand those concepts—or understands them incorrectly or incompletely—will be flawed and less useful than an AI system with a more complete “world model.”<sup>22</sup> And without a broad and diverse array of training data, it will be impossible for AI systems to develop anything close to complete and accurate models of these concepts.

The more limitations we place on the training data available to AI models, the less useful the models will be. An AI system trained on licensed stock photographs, for example, may be capable of identifying an elephant in a *National Geographic* spread, but might be incapable of finding an elephant in a cartoon meme or a child’s drawing because it lacks a complete model of the concept of an “elephant.” A model trained on licensed film scripts might understand the films’ characters, but would be ignorant of the cultural significance those characters: such a model might be able to answer questions about the films—*i.e.* “Who was Darth Vader’s son?”—but would be unaware that the name “Darth Vader” has, since 1977, become a trope for a powerful and malicious political figure.<sup>23</sup> Worse still, a model trained on public domain books will fail to understand modern customs, language, and values—and (even more problematically) will “learn” the discriminatory biases inherent in texts published in the late 19<sup>th</sup> and early 20<sup>th</sup> centuries.<sup>24</sup>

Without access to a broader array of content—including, for example, internet forum comments or memes that use the Darth Vader trope for political or social commentary—a model’s

---

<sup>20</sup> NMPA Comment at 20–21.

<sup>21</sup> See, e.g., Copyright Alliance Comment at 59–60; Getty Comment at 5.

<sup>22</sup> See, e.g., Meta, “Yann LeCun on a Vision to Make AI Systems Learn and Reason like Animals and Humans, Meta Blog,” Meta Blog (Feb. 23, 2022), available at: <https://ai.meta.com/blog/yann-lecun-advances-in-ai-research/> (AI systems need “internal models of how the world works”).

<sup>23</sup> See, e.g., Alex Luhn, “The ‘Darth Vader’ of Russia: Meet Igor Sechin, Putin’s Right-Hand Man,” Vox.com (Feb. 8, 2017), available at: <https://www.vox.com/world/2017/2/8/14539800/igor-sechin-putin-trump-sanctions-oil-rosneft-tillerson-secretary-of-state-kremlin>.

<sup>24</sup> Cf. The White House, Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence § 1 (Oct. 30, 2023) (stressing the importance of preventing “discrimination” and “bias” in AI systems); see also Amanda Levendowski, How Copyright Law Can Fix Artificial Intelligence’s Implicit Bias Problem, 93 Wash. L. Rev. 579, 615 (2018) (“Most public domain works were published prior to 1923, back when the ‘literary canon’ was wealthier, whiter, and more Western than it is today. . . . A dataset reliant on works published before 1923 would reflect the biases of that time, as would any AI system trained with using that dataset.”).

understanding of these concepts will be incomplete. That incompleteness, in turn, will constrain our ability to build AI systems with which humans can converse and collaborate. And, as Meta and others already explained, because much of the content required to illuminate these semantic nuances is user-generated or orphaned, it is virtually impossible to license.<sup>25</sup>

Somewhat ironically, limiting the amount of training data from which models can learn also increases the likelihood that the models will simply regurgitate features of the training data itself. A model whose training data includes only a small subset of modern fiction—licensed, for example, by Bloomsbury Publishing, the publisher of the *Harry Potter* series—might be more likely to assume that all “wizards” are trained at British boarding schools. By contrast, a model trained on a broader scope of content might have a more complete and holistic understanding of the word “wizard”—informed not only by J.K. Rowling’s works, but also by the hundreds of other fictional wizards that enrich any human’s understanding of that word. When asked to create a tale about a wizard, the former model would most likely create a story that simply regurgitates narrative elements from *Harry Potter*. The latter model, by contrast, is far more likely to employ its more robust understanding of what a “wizard” can be to create something entirely novel.

In this sense, AI models are similar to search engines. A search engine capable of searching only a limited subset of the internet is significantly less useful than a search engine that compiled its index by crawling all (or substantially all) of the internet’s publicly accessible websites. The very function of the tool requires a near-universal scope. Forcing the search engine to secure a license from each and every rightsholder with a copyright interest in content on a public website would, in effect, deprive the public of this important mechanism for navigating the web.<sup>26</sup> So too for AI. The AI systems that will be central to our economy in the coming decades will be those that possess a complete and accurate model of the concepts we use to talk to each other and share ideas. Models trained on a small subset of data—*i.e.* content in the public domain or content that AI developers manage to license from rightsholders—will necessarily possess a less complete model of the world and, in turn, be much less effective at collaborating with human users.

For that reason, placing limitations and licensing requirements on AI developers will necessarily result in weaker, less effective AI models. That, in turn, will imperil the United States’ current role as the home for AI innovation, opening the door for other countries—many of which are currently considering how to encourage innovation and earn a leading role in this space—to become the preferred destination for AI development.

#### B. There is No Evidence that a Licensing Market is Feasible

Second, there is no evidence that licensed or public domain data is sufficient to build a useful state-of-the-art Generative AI model capable of competing with available alternatives. For example, Adobe—the creator of the Firefly model discussed in several comments—itsself explained to the Office that “an AI system trained on a small dataset is at greater risk of producing

---

<sup>25</sup> Meta Initial Comment at 19–20; *see also* TechNet Comment at 9–10.

<sup>26</sup> *Cf. Field v. Google Inc.*, 412 F. Supp. 2d 1106, 1117–23 (D. Nev. 2006) (creation of “snapshot[s]” of webpages for internet search was fair use); *see also Kelly v. Arriba Soft Corp.*, 336 F.3d 881, 817–22 (9th Cir. 2003) (creation of copies of images created for search engine was fair use).



wrong or unsatisfactory results or reproducing harmful biases that exist within the dataset.”<sup>27</sup> The existence of a handful of Generative AI models trained on licensed or public domain data is not evidence that all Generative AI models may be developed the same way, or that models trained on limited subsets of data will perform as well as models trained on a broader array of examples.

Ultimately, whether it is possible to train a competent Generative AI model using only public domain or licensed data will depend on a number of fact-specific considerations, including the medium of the model’s output. Several commenters, for example, pointed to a number of music-generation models as evidence that licensing is possible, including Meta Platforms’ MusicGen and Stability AI’s Stable Audio, both of which were trained on licensed data.<sup>28</sup> Meta created MusicGen using only 20,000 hours of music.<sup>29</sup> The resulting model, while impressive, is a simple and largely experimental tool that recognizes key terms in a user prompt and then synthesizes generic musical samples to match the prompt’s description.<sup>30</sup> It hardly constitutes an all-purpose music generation tool, and it is an even farther cry from the general-purpose, “intelligent” models that will power our AI industry in the coming decades.

A large language model, by contrast, faces the much more challenging task of capturing the meanings of hundreds of thousands of words, the complex grammatical rules of hundreds of different languages, and the virtually infinite universe of ideas and concepts that can be expressed by human language. (This most likely explains why none of the comments submitted to the Office suggest that there exist any general-purpose large language models that have been trained on only public domain or licensed data.) And multi-modal models—*i.e.* those that can interpret and output content in any medium—will require an even broader scope and diversity of training data.

Moreover, the fact that some AI developers “might be willing to purchase licenses in order to engage in this transformative use” is “irrelevant” because “[l]ost licensing revenue counts under Factor Four only when the use serves as a substitute for the original.”<sup>31</sup> AI model training does not substitute for the original work, but rather contributes to a new body of works independent from the original. As Meta explained in its opening comments, this is a quintessential transformative fair use,<sup>32</sup> and it is well established that rightsholders “may not preempt exploitation of transformative markets . . . by actually developing or licensing others to develop those markets.”<sup>33</sup> Fair use is a critical feature of our copyright law that “help[s] to keep a copyright

---

<sup>27</sup> Adobe Comment at 2.

<sup>28</sup> See NMPA Comment at 20–21; News/Media Alliance Comment at 23.

<sup>29</sup> See Meta, “Simple and Controllable Music Generation,” (Nov. 7, 2023), available at: <https://arxiv.org/abs/2306.05284>. Similarly, Stability AI’s Stable Audio model was trained using roughly 19,500 hours of audio. See “Stability.AI, Stable Audio: Fast Timing-Conditioned Latent Audio Diffusion,” (Sept. 13, 2023), available at: <https://stability.ai/research/stable-audio-efficient-timing-latent-diffusion>.

<sup>30</sup> See Meta, “MusicGen: Simple and Controllable Music Generation”, available at: <https://ai.honu.io/papers/musicgen/>.

<sup>31</sup> *Authors Guild, Inc. v. HathiTrust*, 755 F.3d 87, 100 (2d Cir. 2014).

<sup>32</sup> Meta Initial Comment at 12–14.

<sup>33</sup> *Castle Rock Entm’t, Inc. v. Carol Pub. Grp., Inc.*, 150 F.3d 132, 145 n.11 (2d Cir. 1998).

monopoly within its lawful bounds.”<sup>34</sup> The existence of a handful of narrow licensing deals in the general domain of Generative AI is not a sufficient basis to narrow that important doctrine and expand the scope of copyright’s traditional monopoly.

In any case, like a fair use analysis, any determination as to the feasibility of licensing content for the development of AI models—or, for that matter, as to the amount and diversity of content necessary to train a useful model—“calls for case-by-case analysis.”<sup>35</sup> Similarly, opining on the general merits of the fair use defense in this context makes little sense, particularly due to the substantial variation between different kinds of model types—from basic music generation models (which may be able to extract general attributes about musical genres from a relatively small set of samples), to highly complex large language models (which seek to build “world models” that capture a broad and complex array of concepts). For that reason, it would be inappropriate for the Office to opine on these issues in the abstract, particularly when it lacks the evidence to do so.

#### **IV. The “Quality” of Any Individual Piece of Training Data is Immaterial**

Several commenters claimed that their works are of particularly high “quality,” and suggested that this somehow affects the fair use analysis.<sup>36</sup> To be sure, it is well recognized in the field of data science that the machine learning models do not work without “high quality” data—meaning that the data must be complete, de-duplicated, and free of errors.<sup>37</sup> Put simply, a weather prediction model trained on a complete and reliable dataset will be more effective than a weather prediction model trained on a dataset that is incomplete and riddled with inaccurate and inconsistent meter readings.<sup>38</sup> Similarly, a Generative AI system like a large language model trained on a dataset containing sentences with proper grammar, vocabulary, and syntax will be more useful than a language model trained on garbled or incomplete text. But that is not the same definition of “quality” as used by the commenters.

Instead, commenters’ arguments regarding data “quality” appear to suggest that copyrighted works with more aesthetic or commercial appeal—like novels or popular works of visual art—are necessarily more valuable as training data. To be sure, such works might feature

---

<sup>34</sup> *Google LLC v. Oracle Am., Inc.*, 141 S. Ct. 1183, 1198 (2021); *see also Eldred v. Ashcroft*, 537 U.S. 186, 219–21 (2003) (noting importance of fair use defense to the “traditional contours of copyright protection”).

<sup>35</sup> *Campbell v. Acuff-Rose Music, Inc.*, 510 U.S. 569, 577 (1994).

<sup>36</sup> *See, e.g.*, News Corp Comment at 3 (noting that “[j]ournalistic works are [] exceptionally well-written” and “thoughtfully conceived”).

<sup>37</sup> *See, e.g.*, Lukas Budach, et al., “The Effects of Data Quality on Machine Learning Performance,” (Nov. 9, 2022), available at: <https://arxiv.org/pdf/2207.14529.pdf>; *see also* Lora Aroyo, et al., “Data Excellence for AI: Why Should You Care,” (Nov. 19, 2021) <https://arxiv.org/ftp/arxiv/papers/2111/2111.10391.pdf>. Some commenters cite these and similar scientific studies in support of their argument that their content is more “valu[able]” to Generative AI models than other content. *See, e.g.*, News Corp Comment at 2 & n.3. But these studies address data “quality” from a data science perspective, and discuss the importance of attributes like completeness and de-duplication to machine learning in general. *See supra*. None suggest that content with more commercial appeal is more valuable for AI training purposes.

<sup>38</sup> *Cf.* Meta Initial Comment at 2–3.

richer narratives, character development, imagery, or diction. But Generative AI models do not exploit the expressive content of the works included in their training data, and do not in any way capture the expressive features that make some copyrighted works more aesthetically or commercially appealing than others. Instead, Generative AI models use training data to identify and learn from patterns gleaned from a broad spectrum of data. As a result, training data is *substitutable*: the exclusion of any specific work from the training corpus will have a negligible effect on the model’s ultimate function, as long as the model has other examples to learn from.

For example, a model whose training data included the sentence “All that is gold does not glitter; Not all those who wander are lost”<sup>39</sup> might learn the relationship between the words “gold” and “glitter” after seeing those words appear in thousands of other sentences. It may even begin to understand the figurative, metaphorical meaning of the word “gold” after seeing that meaning demonstrated in other texts, like restaurant or film reviews. But J.R.R. Tolkien’s evocative diction and sentence structure is not of a higher “quality” to a Generative AI model than less expressive sentences like “Gold glitters when exposed to light” or “This restaurant’s sticky toffee pudding is pure gold.”

Put simply, that a particular work might be “exceptionally well-written,” “thoughtfully conceived,” or deserving of critical acclaim says nothing about that work’s value as training data.<sup>40</sup> For that reason, a work’s commercial or aesthetic value is irrelevant to the question whether using it to train a Generative AI model qualifies as a fair use. As the Second Circuit explained in *Google Books*, a “secondary use” that merely extracts “information about the original, rather than replicating protected expression” is a fair use, regardless of the expressive content of the original.<sup>41</sup>

## V. Generative AI Furthers the Constitutional Purpose of Copyright

Virtually all commenters who oppose application of fair use in this context overlook that the reason fair use exists is to safeguard “copyright’s very purpose, ‘[t]o promote the Progress of Science and useful Arts.’”<sup>42</sup> Courts for decades have used fair use to “avoid rigid application of the copyright statute when . . . it would stifle the very creativity which that law is designed to foster.”<sup>43</sup> Using copyright law to block or hinder the development of Generative AI is impossible to reconcile with this principle. The goal and function of Generative AI models is to create new, different content by “build[ing] upon” the body of human knowledge that currently exists.<sup>44</sup> That overall purpose is entirely consistent with the “Progress of Science and useful Arts.”<sup>45</sup> Indeed, creation of new and different works is the singular and primary goal of copyright law.<sup>46</sup> Moreover,

---

<sup>39</sup> J.R.R. Tolkien, *The Fellowship of the Ring* (1954).

<sup>40</sup> News Corp Comment at 3.

<sup>41</sup> *Authors Guild v. Google, Inc.*, 804 F.3d 202, 220 (2d Cir. 2015).

<sup>42</sup> *Campbell v. Acuff-Rose Music, Inc.*, 510 U.S. 569, 575 (1994) (quoting U.S. Const. Art. 1, § 8, cl. 8).

<sup>43</sup> *Stewart v. Abend*, 495 U.S. 207, 236 (1990).

<sup>44</sup> *Campbell*, 510 U.S. at 575.

<sup>45</sup> U.S. CONST. Art. 1, § 8, cl. 8.

<sup>46</sup> *Sony Corp. of Am. V. Universal City Studios, Inc.*, 464 U.S. 417, 429 (1984) (“reward to the owner” is “a secondary consideration”).



Generative AI will not only unlock a new age in human creative production—it will also vastly advance human scientific endeavors, including by saving lives through advances in medicine and healthcare<sup>47</sup> and solving existential global issues like climate change.

The only reason Generative AI implicates the Copyright Act at all is that, unlike human learning, AI training requires the creation of intermediate copies of the content from which the models learn. Those intermediate copies do not in any way prejudice rightsholders’ ability to “secure a fair return” on their labors by selling their works in the marketplace, as rightsholders have done for centuries.<sup>48</sup> Indeed, these intermediate copies are often temporary, retained only as long as necessary to train the models. The argument that the vast societal benefits Generative AI will bring must take a backseat simply because the process of training requires the creation of intermediate copies that never see the light of day is precisely the kind of “rigid application of the copyright statute” that fair use exists to avoid.<sup>49</sup> Copyright, in other words, is not, and has never been, “an inevitable, divine, or natural right that confers on authors the absolute ownership of their creations.”<sup>50</sup> Copyright owners have never had the right to block others from using their works to derive the basic tenets of human knowledge necessary to advance “the Progress of Science and useful Arts.”<sup>51</sup> To use copyright law to block or impede the development of this new, welfare-enhancing technology would turn copyright upside-down.

---

<sup>47</sup> Scientists and researchers continue to develop innovative AI healthcare tools using Meta’s open-source LLMs that promise to have significant real-world benefits. For example, Radiology-Llama2 is an LLM specialized for radiology applications being developed by researchers at a number of institutions including Harvard University, Massachusetts General Hospital, and the Mayo Clinic. See Zhengliang Liu et al., “Radiology-Llama2: Best in Class Large Language Model for Radiology,” (Aug. 29, 2023), available at: <https://arxiv.org/pdf/2309.06419.pdf>.

<sup>48</sup> *Twentieth Century Music Corp. v. Aiken*, 422 U.S. 151, 156 (1975); see also *Campbell*, 510 U.S. at 576 (asking whether the use “prejudice[s] the sale, or diminish[es] the profits, or supersede[s] the objects, of the original work”).

<sup>49</sup> *Stewart*, 495 U.S. at 236.

<sup>50</sup> *Cambridge Univ. Press v. Patton*, 769 F.3d 1232, 1256 (11th Cir. 2014) (cleaned up).

<sup>51</sup> U.S. CONST. Art. 1, § 8, cl. 8.