



SERVING AUTHORS WORLDWIDE
AU SERVICE DES AUTEURS DANS LE MONDE
AL SERVICIO DE LOS AUTORES EN EL MUNDO

SG23-0959

Source language: English / Written on: 01/12/2023

Document prepared by CISAC Secretariat (CISAC)

CISAC Reply Comment to the US Copyright Office's Notice of Inquiry (NOI) on Copyright and Artificial Intelligence

The International Confederation of Societies of Authors and Composers (“CISAC”) submits this reply comment to react and respond to points raised during the first round of public comments made for the US Copyright Office’s Notice of Inquiry (“NOI”) on Copyright and Artificial Intelligence (“AI”). The purpose of this reply is to assist the USCO in reaching the most comprehensive understanding of the issues generated at the intersection of AI and Copyright, in the interest of all stakeholders that are or may be affected by the rise of AI technologies in society.

CISAC is the leading worldwide organisation of authors’ societies, representing more than 5 million creators from all geographic areas and all artistic repertoires (including music, audiovisual, drama, literature, and visual arts) through our 225 member organisations. The diversity of our membership, along with our longstanding history of safeguarding the interests of creators internationally, permits us to advocate the interests of a significant number of affected parties in the search for consensus on pressing copyright and authors’ rights issues. We provide below five statements, which may serve as guideposts for the USCO in its current inquiry.

1. Unregulated or under-regulated AI technologies pose significant risks to creators.

Based on preliminary observations drawn from the submissions of AI developers and large companies investing in AI, the risks of unregulated development of AI technologies have been grossly understated. Though copyright issues were the main focus of this NOI, Microsoft, Google, Meta, Stability AI, OpenAI,¹ and more, attested to the benefits of AI in the fields of healthcare, finance, and scientific research, along with any other generally beneficial potential uses of AI, such as improving productivity and promoting language learning.² While AI is certainly promising in those fields, we find that the dangers of underregulating the development of AI were not fully represented.³

Most importantly, the lists of AI’s societal benefits obscure the issue that is the main subject of the NOI: that the use of AI and its effect on the creative practices of most of today’s creators – visual artists, musicians, producers, and more – will be profound, and will generate significant and potentially irreversible negative consequences for the creative marketplace.

¹ Comment of Microsoft, p. 4-5; Google, p. 5-6, Meta p. 9-10; Stability.ai, p. 6; OpenAI, p. 4.

² Comment of OpenAI, p. 4.

³ For example, the case of the use of AI in law enforcement has demonstrated that bias can be introduced in the design of enforcement tools, already having caused significant damage to marginalized communities. See, e.g., Hao, K. (2019), “AI is sending people to jail—and getting it wrong: Using historical data to train risk assessment tools could mean that machines are copying the mistakes of the past.” MIT Technology Review, 21 January 2019. <https://www.technologyreview.com/2019/01/21/137783/algorithms-criminal-justice-ai/#:~:text=Now%20populations%20that%20have%20historically%20been%20disproportionately%20targeted,more%20bias-tainted%20data%20to%20feed%20a%20vicious%20cycle>.

Indeed, the “AI effect” carries with it numerous risks of damaging the livelihoods of creators. Artists relying on their unique repertoires of work to establish their careers are particularly at stake; there is an alarming ease in the process of an average user instructing an AI software to replicate the characteristics, style, and elements of existing works. At present, artists are also unaware of whether or not their works are being used to train AI models, as the process of gathering data on their works is often not disclosed, collected through automated processes such as web scraping and database extraction.

In response, AI developers increasingly rely on technical arguments to justify their indiscriminate use of copyrighted content for training purposes. One argument maintains that copying, *per se*, does not occur because the AI is trained using unprotectable information and relational data extracted from existing works. Others argue that the general style and character of a work are not recognizable as protectable subject matter, and therefore generative output that mimics the style of works used to train an AI system cannot be the grounds for claiming infringement. As elaborated below, both arguments are short-sighted and rely on technicalities rather than accepted policy, particularly regarding the interpretation of the applicability of the fair use exception.

2. The “Fair Use” defense does not apply *de facto* to the use of copyrighted content for AI training purposes, as an analysis of the factors may favor a finding of infringement.

The fair use defense which many AI developers rely on should be applied as a measure in equity, as many exceptions to copyright law usually are.⁴ While we do not attempt here to engage in a comprehensive analysis of the fair use factors (we refer to the comments of ASCAP and Authors Guild for a comprehensive fair use analysis, which we fully support), it is necessary to highlight key arguments providing for the conclusion that, in the specific case of the use of copyrighted materials for training purposes, the fair use test should favor the creators whose works are being used.

Preliminarily, any consideration of the application of an exception to copyright law should be assessed on a fact-specific basis; categorical assumptions that certain types of uses are considered fair by law is not reflected in the U.S. jurisprudence. Academics agree that “...it is impossible to say categorically that inputs and outputs of Generative AI will always be fair use.”⁵ It is therefore important that any discussion of the application (or not) of a fair use exception should rely on factual inquiry rather than mere conjecture.⁶

The initial question is whether an act of copying is involved in the use of copyrighted materials for training foundation models. On one side, the majority of AI developers and companies investing in AI seem to argue that the use of protected content for training purposes does not involve the creation of copies, either by denying outright that a copy is made or that, on a technical level, training a foundation model does not involve copying a work but rather collecting the underlying organization, patterns, relationships, and other similar data points contained within a work, thus falling outside the scope of protectable subject matter under copyright.⁷

The latter, “technical” argument provided by many AI developers is that machine learning processes involve copying mere unprotectable “information” from the source material (i.e., the data underlying the work, such as the relationship between colors in a digital image, contextual use of words, etc.), rather than copying the

⁴ According to international treaties (i.e., TRIPs, Berne Convention), exceptions must further comply with the “three-step test”, namely that a reproduction of a work may be permitted in certain special cases, provided that such reproduction does not conflict with a normal exploitation of the work and does not unreasonably prejudice the legitimate interests of the author.

⁵ Comment of Samuelson, et. al., p. 9.

⁶ We do not attempt here to engage in a complete analysis of fair use, but instead try to provide useful analysis of factual and conceptual issues arising in the determination of whether the use of copyrighted content for AI training purposes can be considered “fair” under the U.S.’s statutory fair use test.

⁷ Problematically, some AI developers have deemed this process “knowledge-gathering” or “learning” as a human would, which may create further implicit bias favoring these automated processes. See, e.g., Grimmelman, J. (2015), “Copyright for Literate Robots” 101 Iowa L. Rev. 657. (“Almost by accident, copyright law has concluded that it is for humans only: reading performed by computers doesn’t count as infringement”).

work itself, in order to train the model. As explained by OpenAI, for instance, in training its large language models (LLMs), “...the process begins by breaking text down into roughly word-length “tokens,” which are then converted into numbers. The model then calculates each token’s proximity to other tokens in the training data—essentially, how near one word appears in relation to any other word.”⁸ A similar process is explained very generally in the submissions of Microsoft,⁹ Google,¹⁰ and Meta.¹¹

However, evidence was submitted during the first round of comments to suggest that unauthorized copying of works may indeed occur in the process of accumulating training data for foundation models. For example, StabilityAI, like many other AI developers, uses publicly available databases for sourcing a large quantity of high-quality image to train its models, namely Stable Diffusion.¹² One such database used by StabilityAI is LAION-5B, “an open dataset of publicly-accessible image links and captions compiled by the European nonprofit Large AI Open Network (LAION).”¹³ According to the comment submitted by Michael Frank, researcher and owner of the website “<https://whatsinlaion.com>”, LAION only distributes the web addresses of images, not the actual images themselves, due in part to the large size of the files, but also perhaps in an attempt to limit its liability.¹⁴ Frank concludes that,

“...machine learning developers utilizing the LAION datasets (such as StabilityAI) must download the images directly from the hosts before beginning training. When downloading images using these indexes, copies are necessarily made (however transient) to facilitate the AI training process. These copies are used without approval of the original rights holder, or the approval of the web provider”.¹⁵

Furthermore, in the AI developers’ selection and filtering of training data, expressive works are sought out specifically for their expressive value,¹⁶ meaning that an original work is capable of being replicated by AI in an identical way to the original expression.¹⁷ According to well-established precedent, the fair use exception has never covered uses of creative content for expressive purposes, as this likely fails the first factor of the statutory test.¹⁸ Likewise, unlike precedent cases relying on the transformativeness of the new use (i.e., the creation of a “search” function through the book digitalization process), the outputs of generative AI are similar-in-kind to the works they were trained on. On this basis, the *use of generative AI in general* cannot be considered transformative according to the standard established in cases like *Authors Guild* and *HathiTrust*.

In addition, while general aspects of style and character are not themselves copyrightable, AI can create, as the final output, near perfect copies of works or deceptively similar works so that the average person cannot discern a difference in the marketplace. The presence of these AI-generated or AI-assisted works on the marketplace creates direct competition with the creators whose works have been used to train the model, which would also likely fail the fair use test’s fourth factor.¹⁹ As mentioned in the previous submission,

⁸ Comment of OpenAI, p. 5.

⁹ Comment of Microsoft p. 6.

¹⁰ Comment of Google, p. 9.

¹¹ Comment of Meta, p. 2.

¹² Comment of StabilityAI, p. 10.

¹³ Id.

¹⁴ Comment of Michael Frank, “Use of Copyrighted Material in the Training of Stable Diffusion”, p. 3.

¹⁵ Id.

¹⁶ Comment of Michael Frank, p. 4-6 (commenting on the presence of bias for aesthetic quality in selecting data for training models, citing the use of LAION’s “Aesthetic Scoring” image ranking system, which adds values to images of between 1 to 10 for models to rank the images for visual quality, with 10 representing the most aesthetically-pleasing images. The result was that “images that fell below a rank of 4.5 were not included in the creation of aesthetic datasets.”).

¹⁷ An amended class-action lawsuit was filed on 29 November by a group of visual artists against Stability AI, Midjourney and other companies for allegedly misusing their work to train generative artificial intelligence systems. From the complaint, “AI image products are primarily valued as copyright-launders devices, promising customers the benefits of art without the costs of artists.” Full text of amended complaint provided by Reuters:

<https://fingfx.thomsonreuters.com/gfx/legaldocs/znbnzrgyzpl/AI%20COPYRIGHT%20LAWSUIT%20amended.pdf>.

¹⁸ See *Authors Guild v. Google, Inc.*, 804 F.3d 202 (2d Cir. 2015); *Authors Guild, Inc. v. HathiTrust*, 755 F.3d 87 (2d Cir. 2014); and *AV ex rel. Vanderhye v. iParadigms, LLC*, A.V. ex rel. Vanderhye v. iParadigms, LLC, 562 F.3d 630 (4th Cir. 2009) (representing cases of fair use which did not involve exploitation of works for their creative content.) See also, Comment of Authors Guild.

¹⁹ See also, Comments of ASCAP; Authors Guild; Copyright Alliance.

generative AI outputs may become available at much lower costs and in higher quantities, driving down the overall cost of creative labor. This would potentially create a substitutive effect of generative AI outputs over human-created original works on the marketplace, as consumers may choose such works over the original copyrighted content, directly impacting the market for the original work.

Finally, the third factor regarding the amount and substantiality of the use of copyrighted works with respect to their use as training data for foundation models, commented on in CISAC's first submission²⁰, should weigh in favor of copyright rightsholders due to the fact that the data ultimately selected for training foundation models necessarily requires the use of works in their entirety – otherwise, the relational data and other information gathered for the purposes of training the model would be incomplete.

To sum up the results of this brief fair use discussion and analysis, we conclude that the fair use defense does not absolve AI developers of any claims of copyright infringement (or failure to license copyrighted content), due to the uncertainty of its application to the use of copyright protected works for training foundation models. Therefore, AI developers shall take a fair, transparent, and responsible approach and engage with rightsholders to negotiate licensing agreements for the use of copyright protected content.

3. TDM exceptions as they are implemented internationally do not permit *de facto* the use of copyrighted content for AI training purposes.

As a general observation, based on the reply submissions, AI developers lean heavily on blanket assumptions that they are able to benefit from the exception of fair use and any other similar copyright exceptions as they appear in national law, but fail to provide specific reasoning or specific factual bases for applying such exceptions. In other words, AI developers fail to acknowledge that “exceptions,” by their very nature, are exceptionally applied.

One particularly significant overgeneralization advanced by AI developers was that the presence of TDM exceptions in the national laws of certain jurisdictions *de facto* permit the use of copyrighted works for training AI. This is an incorrect interpretation of law regarding TDM exceptions and their application. According to the submission of Microsoft, the company mentions that, “[w]e welcome the steps that many jurisdictions are taking, under their national IP frameworks, to clarify that copyright law continues to permit AI training”, citing the jurisdictions of Japan, Singapore, Korea, the European Union and Israel.²¹

However, the language, interpretation and application of TDM exceptions vary greatly between these jurisdictions. Additionally, it is unclear to what extent these provisions are compatible with international treaties, namely TRIPS and the Berne Convention's three-step test, which applies to exceptions and limitations to the right of reproduction.²² It is therefore not possible to conclude that TDM exceptions provide a blanket permission to use copyrighted content for AI training purposes.

Therefore, even considering the potential application of an existing exception or limitation to copyright, we contend that it is in the best interest of stakeholders that the use of copyrighted materials as training data for foundational models should be licensed in all cases.

4. Voluntary licensing solutions can be developed and implemented successfully.

During the first round of comments, many companies having significant investments in AI voiced concerns regarding the feasibility of licensing copyrighted content for training purposes.²³ Such licensing solutions,

²⁰ Comment of CISAC, p. 3.

²¹ Comment of Microsoft, p. 7.

²² Comment of Gervais, D., p. 2.

²³ Comments of Microsoft, p. 9; OpenAI, p. 13; Meta, p. 14; Hugging Face, p. 11.

however, can be developed to meet the very specific needs of this market, and ample precedent exists for this conclusion.

In his testimony before Congress, former General Counsel of the US Copyright Office John Baumgarten analogized the current changes introduced into the copyright licensing landscape to developments that occurred at the beginning of the age of reprography.²⁴ In his comparison, he draws attention to the fact that collective licensing solutions had been developed to address all concerns raised by industry representatives who were previously reluctant that licenses could exist in the new age of large-scale industrial copying. He further maintains that, in the context of constructing similar licensing solutions to adapt to current challenges with the use of protected materials for training purposes, “...collective licensing not only assures compensation (including monetization of micro-transactions and flow of funds across borders) to individual creators for the widespread use of their creative works, it also offers great advantages to *users* of copyrighted works....[i]ndeed, collective management infrastructures already exist which can collect and distribute royalties, even in the face of new and complex licensing terms and limitations.”²⁵

As for proof that a market exists for licensing creative content for training purposes in the first place, OpenAI mentions that its models are pre-trained using “nonpublic information that we obtain from third parties through commercial arrangements”, indicating clearly that a market for licensing content for training purposes has already emerged, and that AI developers are willing to negotiate and pay for access to such non-publicly available content.²⁶

CISAC is particularly well-placed to comment on the feasibility of developing new licensing and rights management solutions due to its long history supporting collective management and performing rights organisations in developing universally accepted standards²⁷ for accurately tracing and tracking royalties owed to creators. As supported by its member societies, collective licensing of creative works in the context of their use to train AI is feasible and may provide new opportunities for collaboration and cooperation between all involved stakeholders. For example, OpenAI concluded a license with the Associated Press in July for the use of its news archive in the development of generative AI models.²⁸ Such licensing efforts need to be widely adopted to ensure a sustainable future for creators and innovators alike, and can be effectively fostered through existing infrastructures used for managing rights which are already present across all creative sectors.

5. Assuring the future of human creators also assures the future creative potential of AI, and societies’ future access to creative works.

As mentioned in our previous submission, we maintain that a well-functioning copyright law is capable of balancing private and public interests; when rights, responsibilities and obligations are clearly established, all stakeholders may benefit from such a system.

While it is true that AI has the potential to augment human creative efforts, it should be acknowledged that it can also supplant them. Moving forward, the dangers of unregulated or underregulated AI development can be addressed through cultivating a culture around responsible innovation. This can be accomplished in

²⁴ Baumgarten, J. “Comment during Hearing on Artificial Intelligence and Intellectual Property: Part I — Interoperability of AI and Copyright Law of the Subcommittee on Courts, Intellectual Property, and The Internet of The House of Representatives Committee on the Judiciary (May 17, 2023)” Copyright Alliance. <https://copyrightalliance.org/wp-content/uploads/2023/05/Final-Baumgarten-AI-Letter-to-IP-Subcomm.pdf>.

²⁵ Id.

²⁶ Comment of OpenAI, p. 5.

²⁷ For example, International Standard Musical Work Code (ISWC) and International Standard Recording Code (ISRC), Interested Party Information (IPI), among numerous other technical standards currently accepted industry-wide. CISAC, “Information Services”, <https://www.cisac.org/services/information-services>.

²⁸ AP News, “ChatGPT-maker OpenAI signs deal with AP to license news stories” <https://apnews.com/article/openai-chatgpt-associated-press-ap-f86f84c5bcc2f3b98074b38521f5f75a>.

many ways, but at least requires coordinated efforts and meaningful collaborations between AI developers, creators, governments, and the public to find the most workable solutions for all.

The prevalence of automation in many industries has driven down costs of labor, but workers have not captured any of its benefits: instead, greater margins are kept by companies and are not passed on. However, at this early stage in the development and integration of AI in society, the same outcome is avoidable – when it comes to the use of AI technologies making creative processes cheaper and easier, the benefits of reducing such creative labor costs should not be captured and kept by AI companies but should rather help remunerate the human creators at the source of the AI's training process.

We firmly believe that a more robust creative marketplace for human creators can coincide with, rather than challenge, innovation in the field of AI. Richer and more diverse data sources are the key to the development of better AI models, and the promise of copyright provides human authors with the ability and incentives to make a living from their works. Once creators are able to fully benefit from uses of their works, the promise of copyright can continue to be kept in an AI-driven era.