

September 11, 2024

Digital Fools' Gold – The Rise (and Potential Fall) of Rightsholders' Efforts to Cash In on the GenAI Gold Rush

By [Lee F. Johnston](#)

The dominance of generative artificial intelligence (“GenAI”) in today’s popular culture is undeniable. Practically every 30 second advertising spot during the Olympics touted the mind-blowing capabilities of yet another GenAI-assisted tool promising to take us to new levels of efficiency and creativity. Business journals report that massive amounts of capital investment are being directed not only to companies who create GenAI-assisted tools, but also to the companies like Nvidia, whose products support the development of these AI-assisted tools.

But what about compensating the individuals who provided the fuel—the massive amounts of data needed to train AI models—without which the “AI revolution” would have never happened? Recent court decisions suggest that the answer to this question will turn on whether the harvesting and use of this training data constitutes “fair use” under the Copyright Act.

Brief Overview of GenAI

“Generative AI” systems, such as the Generative Pretrained Transformer (GPT) and Large Language Model Meta AI (LLaMA) language models and the Stable Diffusion and Midjourney text-to-image models, were built by ingesting massive quantities of text and images from the internet. Today’s GenAI models are machine learning models trained on social media posts, books, articles, photos, digital art, music, software, and more. Rather than simply classifying these diverse inputs and generating metadata about them—as previous generations of machine learning systems have—GenAI models can produce new digital artifacts: new text, new art, new music, and new software.¹

Generally speaking, the GenAI-assisted tools are developed and work to create prompt-generated outputs through two transition processes. The first transition is the *training* process. In training, the system is exposed to a large amount of data relevant for its purpose, known as the *training set*. In many instances, this data is harvested from publicly available Internet websites, including social media sites like LinkedIn, Facebook, and Instagram, through an automated process known as data scraping. Thus far, efforts to legally enjoin data scraping of this publicly available (*i.e.*, non-password protected) data via contract (e.g., breaches of website Terms of Service) or other non-copyright legal theories (e.g., unfair competition, unjust enrichment) have been rejected by courts.²

The purpose of training is to extract metadata: information about patterns and relations between elements of the data in the training set.³ The metadata is embodied in a model: a complex set of parameters and “weights” that mathematically represent the extracted patterns. One example of this metadata-based model is Large Language Models (“LLMs”), which represent mathematical patterns and relations between basic elements of language.⁴

At the other end of the process lies the second transition – *generation*. During generation, the GenAI model uses the metadata to create a new specific information output, hopefully one that matches the user’s needs (as dictated by his or her prompt). A prompt is a data input from the user that initiates a process in which the system, by reference to the metadata in the model, constructs a new information output. Within the field of

expression, the process of GenAI can be applied to a growing array of media—including text, speech, images, videos, and music—resulting in an impressive capability of machine creation in these areas.⁵

Recent GenAI Cases

In the past three years, we have seen an explosion of litigation brought by the individual and collective rightsholders, both on behalf of themselves and putative classes, seeking compensation based on both the use of their creative works to train GenAI models and the outputs generated by GenAI tools. Given the number of cases and similarity of claims, courts have consolidated many of the cases. Some of the most noteworthy of these cases include the following:

- *Alter v. OpenAI* (Southern District of New York) – What started as three separate cases brought by three different author groups has been consolidated into a single action against OpenAI and Microsoft. (This case includes Authors Guild and Basbanes). Plaintiffs alleged that OpenAI and Microsoft are liable for copyright infringement arising from the use of plaintiffs’ works to train defendants’ AI models.
- *OpenAI ChatGPT Litigation* (Northern District of California) – Three plaintiff groups, fiction and nonfiction authors, each filed a complaint in the Northern District of California against OpenAI, alleging copyright infringement, vicarious copyright infringement, DMCA violations⁶ and torts related to OpenAI’s GPT models and ChatGPT service. The consolidated cases include *Tremblay v. OpenAI*, *Silverman v. OpenAI*, and *Chabon v. OpenAI*.
- *Anderson v. Stability AI* (Northern District of California) – Visual artists filed this putative class action, alleging direct and induced copyright infringement, DMCA violations, false endorsement and trade dress claims based on the creation and functionality of Stability AI’s Stable Diffusion and DreamStudio, Midjourney Inc.’s eponymous generative AI tool, and DeviantArt’s DreamUp.
- *Doe v. GitHub, Inc.* (Northern District of California) – Anonymous plaintiffs filed this putative class action against GitHub, Microsoft, and OpenAI, alleging that defendants used plaintiffs’ copyrighted open-source code to create Codex and Copilot. Codex is the OpenAI model that powers GitHub’s AI pair programmer, Copilot. The plaintiffs allege that Copilot does not comply with the Open-Source Software licenses governing plaintiffs’ code that was stored on GitHub.
- *Center for Investigative Reporting v. OpenAI* (Southern District of New York) – The Center for Investigative Reporting, a nonprofit news organization, filed a complaint against OpenAI and Microsoft, alleging that the defendants’ use of news articles in training data directly and indirectly infringed plaintiff’s copyrights in the articles.
- *Daily News v. Microsoft* (Southern District of New York) - On April 30, 2024, eight newspaper publishers sued OpenAI and Microsoft, alleging that defendants “purloin[ed] millions of [Plaintiffs’] copyrighted articles without permission and without payment to fuel the commercialization of” their generative AI, including ChatGPT and Copilot.
- *David Millette v. OpenAI, David Millette v. Google, YouTube Inc., David Millette v. Nvidia Corp.* (Northern District of California) – On August 2, 2024, Mr. Millette filed companion class action lawsuits on behalf of YouTube video users and creators, against OpenAI and Google alleging

that the transcription of YouTube videos and use of those transcripts in the training data for defendants' respective GenAI models constituted unfair business practices under California's UCL statute and unjustly enriched defendants. On August 14, 2024, Mr. Millette filed a similar class action complaint against Nvidia, asserting that Nvidia's Cosmos AI technology was improperly trained on YouTube videos.

- *UMG Recordings, Inc. v. Suno, Inc.* (District of Massachusetts) – On June 24, 2024, UMG Recordings filed a lawsuit on behalf of music recording copyright holders against Suno, a generative AI service that produces digital music files based on user prompts that allegedly “compete” with existing copyrighted songs. UMG seeks willful copyright statutory damages of \$150,000 per copyright violation based on its allegations that Suno copied and trained its AI model on more than a decade's worth of the world's most popular sound recording.

As expected, the defendants in these cases have sought dismissal of the state law business tort claims of unfair competition and unjust enrichment, and thus far, courts have been receptive to their Rule 12(b)(6) motions, finding that these claims are preempted under the Copyright Act.⁷ Additionally, defendants have successfully challenged the sufficiency of rightsholders' claims of vicarious and derivative work copyright infringement where direct infringement cannot be plausibly plead as a result of the fact that the AI-generated output lacks “substantial similarity” to the original work.⁸ Finally, courts have dismissed DMCA Section 1202(b) claims where it cannot be plausibly alleged that the defendants distributed the original works, much less did so by removing the Copyright Management Information (“CMI”) on those original works.⁹

All Roads Lead to the Copyright “Fair Use” Defense

Despite defendants' successful motion practice on state law and indirect copyright infringement claims, rightsholders' direct copyright infringement claims have survived where allegations of “substantial similarity” between the original works and the AI-generated output can be plausibly made. In these instances, plaintiffs have pled that, when prompted, AI models have generated near identical excerpts of their original works.¹⁰ As a result, the case dispositive issue—whether the use of their original works and GenAI outputs constitute “fair use”—must await post-discovery summary judgment motions and/or trial.

The analysis that courts use to determine whether a particular use of a copyrighted work constitutes “fair use” is framed in terms of the four factors of § 107 of the Copyright Act of 1976:

- (1) the purpose and character of the use, including whether such use is of a commercial nature or is for nonprofit educational purposes;
- (2) the nature of the copyrighted work;
- (3) the amount and substantiality of the portion used in relation to the copyrighted work as a whole; and
- (4) the effect of the use upon the potential market for or value of the copyrighted work.¹¹

Classic fair uses include parody, commentary, or criticism. These genres are not fair use merely because they change the underlying work or convey some new meaning or message. Rather, the transformations that these genres make are such that the “new” work poses no risk of expressive substitution to the original.¹²

The Supreme Court's 2023 decision in *Andy Warhol Foundation for Visual Arts v. Goldsmith*¹³ provides a helpful definition of “transformative use.” The Court explained that “whether an allegedly infringing use has a further purpose or different character ... *is a matter of degree*, and the degree of difference must be weighed against other considerations, like commercialism.”¹⁴ The Court further explained:

HAYNES BOONE

Most copying has some further purpose, in the sense that copying is socially useful *ex post*. Many secondary works add something new. That alone does not render such uses fair. Rather, the first factor (which is just one factor in a larger analysis) asks ‘whether *and to what extent*’ the use at issue has a purpose or character different from the original. The larger the difference, the more likely the first factor weighs in favor of fair use. The smaller the difference, the less likely.¹⁵

As the Court observed, technical acts of copying that do not communicate an author’s original expression to a new audience have been held to constitute fair use.¹⁶ Examples of non-expressive uses include copying object code to extract uncopyrightable facts and interoperability keys (“reverse engineering”),¹⁷ an automated process of copying student term papers to compare to other papers for plagiarism detection,¹⁸ copying HTML webpages to make a search engine index,¹⁹ copying printed library books to allow researchers to conduct statistical analyses of the contents of whole collections of books,²⁰ and copying printed library books to create a search engine index.²¹

Looking Ahead – GenAI Companies’ Possible Success in Defeating Direct Copyright Infringement Claims via the Fair Use Defense

Armed with this precedent, GenAI defendants may have the upper hand in establishing that their GenAI-assisted tools are shielded from liability under the Copyright Act’s fair use defense. GenAI models that do not, in their ordinary and routine operation, copy (or produce copies of) the original expression in their training data may be found to constitute non-expressive use. To be clear, a GenAI model might be used to create work that is expressive in a First Amendment sense, but the term “non-expressive use” is meant to track copyright’s idea-expression distinction, not broader notions of free expression. GenAI defendants may be able to point to the prior precedent noted above and the recent *Andy Warhol Foundation* case to successfully assert that as long as the original expression in the training data is not transmitted to a new audience, the copying that took place to assemble the training data for GenAI is just as much a non-expressive use as was found by courts in the reverse engineering, student-paper plagiarism detection, and book-copying/search engine contexts.

Following on that logic, GenAI companies may then persuasively argue that deriving uncopyrightable abstractions and associations from the training data and then using that knowledge to create new digital outputs is not just transformative; it is highly transformative.²² And, like other non-expressive uses, the incidental, intermediate reproductions of copyrighted works that are created through assembling training data for a GenAI model do not undermine the copyright owner’s interest in communicating the same original expression to the public. In other words, there is no interference with rightsholders’ interest because the copyright owner’s expression is not conveyed by virtue of the copying being done in an intermediate technical step in the analytical process of building the training set and its use by the GenAI model.²³

HAYNES BOONE

Ultimately, whether these types of arguments will resonate with courts remains to be seen. Stay tuned!

¹ President Biden’s Executive Order 14110 defines GenAI as “the class of AI models that emulate the structure and characteristics of input data in order to generate synthetic content.” Exec. Order No. 14110, 88 Fed. Reg. 75191 §3(p) (Oct. 30, 2023).

² See *hiQ Labs, Inc. v. LinkedIn Corp.*, 31 F.4th 1180, 1202 (9th Cir. 2022) (“Giving companies ... free rein to decide, on any basis, who can collect and use data—data that the companies do not own, that they otherwise make publicly available to viewers, and that the companies themselves collect and use—risks the possible creation of information monopolies that would disserve the public interest.”); see also *Meta Platforms, Inc. v. Bright Data Ltd.*, Case No. 23-cv-00077-EMC, 2024 WL 251406 (N.D. Cal. Jan. 23, 2024) (granting defendant data scraper’s motion for summary judgment on Meta’s breach of contract claim).

³ See Bracha, “Generating Derivatives: AI and Copyright’s Most Troublesome Right,” 25 N.C.J. L & Tech 345, 352-354 (April 2024).

⁴ *Id.*

⁵ *Id.*

⁶ The Digital Millennium Copyright Act (“DMCA”) Section 1202(b) prohibits “the removal or alteration of copyright management information (“CMI”),” such as the title, the author, the copyright owner and other identifying information, from a copyrighted work. Plaintiffs allege that GenAI-assisted tools use processes—during the gathering/ingesting of data from original works and the generation of output—which remove CMI in violation of Section 102(b).

⁷ See, e.g., *Kadrey v. Meta Platforms, Inc.*, Case No. 23-cv-03417-VC, 2023 WL 8039640 (N.D. Cal. Nov. 20, 2023); *Anderson v. Stability AI, Ltd.*, 700 F.Supp.3d 853 (N.D. Cal. 2023); *Tremblay v. OpenAI, Inc.*, Case No. 23-cv-03223-AMO, 2024 WL 3640501 (N.D. Cal. July 30, 2024).

⁸ *Id.*

⁹ *Id.*

¹⁰ In these cases, GenAI defendants have challenged the legitimacy of the prompts used by the rightsholders/plaintiffs to generate “substantially similar” outputs. In particular, GenAI defendants assert that rightsholders have essentially “gamed” their prompts by using verbatim text from the underlying works themselves to generate these infringing outputs. See Memorandum of Law in Support of OpenAI Defendants’ Motion to Dismiss filed on June 11, 2024, in *Daily News, LP et al. v. Microsoft Corp., et al.*, Case No. 1:24-cv-03285-SHS, Docket No. 82 (USDC SDNY).

¹¹ 17 U.S.C. § 107.

¹² See *Campbell v. Acuff-Rose Music, Inc.*, 510 U.S. 569, 594 (1994).

¹³ *Andy Warhol Found. for Visual Arts v. Goldsmith*, 143 S. Ct. 1258 (2023).

¹⁴ *Id.* at 1273.

¹⁵ *Id.* at 1275 (emphasis in original) (citations omitted) (citing *Campbell*, 510 U.S. at 579).

¹⁶ *Id.* at 1274.

¹⁷ See *Sega Entertainment, Ltd. V. Accolade, Inc.*, 977 F.2d 1510, 1514 (9th Cir. 1992); see also *Sony Computer Ent., Inc. v. Connectix Corp.*, 203 F.3d 596, 601 (9th Cir. 2000).

¹⁸ See *A.V. ex rel. Vanderye v. iParadigms, LLC*, 562 F.3d 630, 634 (4th Cir. 2009).

¹⁹ See *Perfect 10, Inc. v. Amazon.com, Inc.*, 508 F.3d 1146, 1155-57 (9th Cir. 2007).

²⁰ See *Authors Guild, Inc. v. HathiTrust*, 755 F.3d 87, 105 (2nd Cir. 2014); *Authors Guild, Inc. v. Google, Inc.*, 804 F.3d 202, 209 (2nd Cir. 2015).

²¹ See *HathiTrust*, 755 F.3d at 91.

²² See *A.V. ex rel. Vanderye*, 562 F.3d at 640 (“The district court ... correctly determined that the archiving of plaintiffs’ papers was transformative and favored a finding of ‘fair use.’ iParadigms’ use of these works was completely unrelated to expressive content and was instead aimed at detecting and discouraging plagiarism.”); *HathiTrust*, 755 F.3d at 97 (“[W]e conclude that the creation of a full-text searchable database is quintessentially transformative use.”); *Authors Guild*, 804 F.3d at 216-17 (“We have no difficulty concluding that Google’s making of a digital copy of Plaintiffs’ books for the purpose of enabling a search for identification of books containing a term of interest to the searcher involves a highly transformative purpose, in the sense intended by *Campbell*.”)

²³ See *Sag*, “Fairness and Fair Use in Generative AI,” 92 *Fordham L. Rev.* 1887, 1913-1917 (2024).