

JAMIE ODELL*

Training on Headlines: *The New York Times*, OpenAI, and the Copyright Implications of AI Data Usage

Introduction.....	204
I. The Institutional and Technical Foundations of OpenAI’s ChatGPT.....	206
A. The Development of OpenAI’s Corporate Identity	206
B. The Technology Behind OpenAI’s Generative AI ChatGPT	208
C. The Copyrightability of GPT-4’s Training Dataset.....	211
II. The <i>Times</i> ’s Allegations.....	213
III. OpenAI’s Defenses	215
IV. Merits of the Legal Claims.....	216
A. Direct Copyright Infringement	216
B. Application of the Fair Use Defense.....	218
1. Purpose and Character	219
2. Nature of the Copyrighted Works.....	223
3. Amount and Substantiality.....	225
4. The Use’s Effect on the Market Value of the Copyrighted Work.....	227
V. Broader Legal and Social Consequences	228
A. The Possibility of Algorithmic Gatekeeping	229
B. The Antithesis of the Copyright Clause.....	233

* J.D. Candidate, 2026, University of Oregon School of Law. The author thanks the faculty of the University of Oregon School of Law and the staff of Oregon Law Review for their invaluable feedback and guidance during this process. The author also expresses gratitude to her friends and family; this piece would not be possible without their endless support. The author dedicates this Comment to the creatives of the world—your labor is irreplaceable. *Please note:* At the time this author wrote this Comment, the legal doctrines surrounding copyright law and AI data usage remain unsettled law.

C. More Than Just a © Symbol	237
Conclusion	238

INTRODUCTION

In late November of 2023, OpenAI, Inc. (“OpenAI”) released its chatbot, “ChatGPT,” which quickly went viral as users discovered its ability to perform almost any writing task, including writing stories, coding computer programs, generating brand slogans, and planning for vacations.¹ Within five days, ChatGPT had over one million users.² As of April 2025, OpenAI had upgraded its Generated Pre-Trained Transformer (“GPT”) model more than four times,³ and has come under judicial scrutiny for using copyrighted materials to train its generative artificial intelligence models (“GenAI”) without the copyright owner’s license or consent.⁴

Many plaintiffs, including the Authors Guild—represented by notable authors such as George R.R. Martin and John Grisham—Sarah Silverman, *The New York Times*, and others have initiated legal actions challenging artificial intelligence (“AI”) companies’ use of their copyrighted works to train GenAI algorithms without authorization.⁵ Particularly, the *Times* argues that OpenAI has incorporated millions of its articles into GPT training datasets without compensation, credit, or consent.⁶ Furthermore, the *Times* asserts that OpenAI’s models generate outputs that closely resemble its copyrighted content, thereby functioning as a direct competitor.⁷ Among the remedies sought, the *Times* requests that all its copyrighted materials contained in GPT

¹ Bernard Marr, *A Short History of ChatGPT: How We Got to Where We Are Today*, FORBES (May 19, 2023, at 01:14 ET), <https://www.forbes.com/sites/bernardmarr/2023/05/19/a-short-history-of-chatgpt-how-we-got-to-where-we-are-today/> [https://perma.cc/2HSN-YP9A].

² *Id.*

³ *ChatGPT – Release Notes, Sunsetting GPT-4 in ChatGPT*, OPENAI (Apr. 10, 2025), <https://help.openai.com/en/articles/6825453-chatgpt-release-notes> (on file with the Oregon Law Review) (“Effective April 30, 2025, GPT-4 will be retired from ChatGPT and fully replaced by GPT-4o.”).

⁴ See First Consolidated Complaint at 2–3, *Authors Guild v. OpenAI, Inc.*, 345 F.R.D. 585 (S.D.N.Y. 2024) (No. 23-CV-08292); see also Complaint at 1, *Silverman v. OpenAI, Inc.*, No. 23-CV-03416 (N.D. Cal. Jul. 7, 2023); Complaint at 2–3, *New York Times Co. v. Microsoft Corp.*, 777 F. Supp. 3d. 283 (S.D.N.Y. 2025) (No. 23-CV-11195) [hereinafter *New York Times Complaint*].

⁵ First Consolidated Complaint, *Authors Guild*, *supra* note 4; Complaint, *Silverman*, *supra* note 4; *New York Times Complaint*, *supra* note 4.

⁶ *New York Times Complaint*, *supra* note 4, at 59.

⁷ *Id.*

training sets be destroyed.⁸ In response, OpenAI argues its use of copyrighted materials is fair use because the works are used to train large language models (“LLMs”) and are not intended to act as a substitute for the original works.⁹

As artificial intelligence continues to advance, the ongoing legal battle between *The New York Times* and OpenAI highlights the diverging, complex tensions at the heart of copyright law—balancing innovation against the constitutional promise to protect original expression.¹⁰ The Constitution’s Copyright Clause “promote[s] the progress of Sciences and useful Arts, by securing for limited Times to Authors and Inventors the exclusive Right to their respective Writings and Discoveries.”¹¹ The Copyright Clause thus requires a balance between the societal value of protecting tangible expressions of creativity against the need for flexibility in technological innovations.¹²

This Comment discusses the copyright law concerns associated with the use of copyrighted materials within the training datasets for AI. Part I will discuss OpenAI’s shift in corporate identity before defining and explaining the functionality of OpenAI’s GPT models. Part I will also examine the types of copyrighted works involved in AI training and the methods by which these copyrighted materials are gathered and curated. Part II will analyze the legal and factual claims presented in the *Times*’s complaint against defendants OpenAI, Inc. and Microsoft Corporation, emphasizing the principal allegation of direct copyright infringement. Part III will explain OpenAI’s response to the *Times*’s complaint. Part IV will assess the *Times*’s allegation of direct copyright infringement and assess OpenAI’s fair use defense, examining (1) the purpose and character of the use; (2) the nature of the original work; (3) the quantitative and qualitative scope of the works used; and (4) the impact of the use on the copyrightable material’s market value. Finally,

⁸ *Id.* at 68.

⁹ Defendant’s Memorandum in Support of Motion to Dismiss at 2–3, *New York Times Co. v. Microsoft Corp.*, 777 F. Supp. 3d. 283 (S.D.N.Y. 2025) (No. 23-CV-11195) (“There is a genuinely important issue at the heart of this lawsuit—critical not just to OpenAI, but also to countless start-ups and other companies innovating in this space—that is being litigated both here and in over a dozen other cases around the country (including in this Court): whether it is fair use under copyright law to use publicly accessible content to train generative AI models to learn about language, grammar, and syntax, and to understand the facts that constitute humans’ collective knowledge.”).

¹⁰ U.S. CONST. art. I, § 8, cl. 8.

¹¹ *Id.*

¹² *See id.*

Part V will reflect on the enduring significance of this case, ultimately advocating for a verdict favoring the *Times* as the outcome most consistent with the constitutional aims of copyright law.

I

THE INSTITUTIONAL AND TECHNICAL FOUNDATIONS OF OPENAI'S CHATGPT

A. The Development of OpenAI's Corporate Identity

Understanding the copyright implications of training data used in LLMs like OpenAI's GPT models requires a technical and contextual background. Across the past decade, OpenAI has transitioned from being an open-source nonprofit corporation committed to transparency and broad social benefits¹³ to a closed-source for-profit corporation focused on rapid technological advancement.¹⁴ The legal questions surrounding generative AI in the context of copyright law—what data was used, who owns it, how was the data collected—perhaps ties directly into OpenAI's updated priorities.¹⁵

OpenAI was founded in 2015 by a variety of researchers and scientists.¹⁶ Initially, OpenAI was a nonprofit organization, free from any financial obligation to turn a profit, with a goal to advance AI technology to “benefit humanity as a whole.”¹⁷ However, by 2019, OpenAI became a for-profit organization when it added a profit-

¹³ *Introducing OpenAI*, OPENAI (Dec. 11, 2015), <https://openai.com/index/introducing-openai/> (on file with the Oregon Law Review) (OpenAI's original mission statement: “Our goal is to advance digital intelligence in the way that is most likely to benefit humanity as a whole, unconstrained by a need to generate financial return. Since our research is free from financial obligations, we can better focus on a positive human impact.”).

¹⁴ Dan Milmo, *Why Is OpenAI Planning to Become a For-Profit Business and Does It Matter?*, THE GUARDIAN (Sep. 26, 2024, at 13:27 ET), <https://www.theguardian.com/technology/2024/sep/26/why-is-openai-planning-to-become-a-for-profit-business-and-does-it-matter> [<https://perma.cc/5M9T-3R9P>].

¹⁵ See Jack Hardinges, Elena Simperl & Nigel Shadbolt, *We Must Fix the Lack of Transparency Around the Data Used to Train Foundation Models*, HARV. DATA SCI. REV. (May 31, 2024), <https://hdsr.mitpress.mit.edu/pub/xau9dza3/release/2> [<https://perma.cc/CK3Q-QQ2Y>] (“Information about training data is also vital to lawmakers’ attempts to assess whether foundation models have ingested personal data or copyrighted material.”); Tim O’Reilly, *You Can’t Regulate What You Don’t Understand*, O’REILLY (June 15, 2023), <https://www.oreilly.com/radar/you-cant-regulate-what-you-dont-understand/> [<https://perma.cc/W34T-T4PX>]; Chloe Xiang, *OpenAI’s GPT-4 Is Closed Source and Shrouded in Secrecy*, VICE (Mar. 16, 2023, at 09:21 PT), <https://www.vice.com/en/article/openais-gpt-4-is-closed-source-and-shrouded-in-secrecy/> [<https://perma.cc/H3X8-L2VR>].

¹⁶ Karl Montevirgen, *OpenAI*, BRITANNICA MONEY (Aug. 26, 2025), <https://www.britannica.com/money/OpenAI> [<https://perma.cc/4RH3-LPAT>].

¹⁷ *Introducing OpenAI*, *supra* note 13.

making subsidiary.¹⁸ As of January 2025, OpenAI's funding round was projected to raise \$40 billion from investors, potentially resulting in a total business valuation of up to \$340 billion.¹⁹ This significant increase in investment highlights the company's financial incentives to expand its for-profit activities notwithstanding its nonprofit origins. This corporate refinancing comes nearly a year after the departure of several top executives and the dissolution of OpenAI's Superalignment team—a team previously dedicated to mitigating long-term AI-related risks; namely, “security, monitoring, preparedness, safety and societal impact.”²⁰ These developments suggest a prioritization of rapid AI advancement over cautious governance, raising questions about the implications of OpenAI's evolving corporate identity on the broader landscape of AI safety and ethics.

Moreover, in the beginning stages, OpenAI adopted an open-source strategy, as opposed to a closed-source strategy. The open-source strategy reflected its initial mission statement, which promised to share its research, code, and potential patents with the world.²¹ In the context of AI development, an open-source AI system is one that provides information, like training data, source code, and weights of parameters, “to the extent that ‘a skilled person can recreate a substantially equivalent system using the same or similar data.’”²² Conversely, a closed-source AI system is one that is closely and strictly controlled by

¹⁸ Milmo, *supra* note 14.

¹⁹ Hayden Field & Kate Rooney, *OpenAI in Talks to Raise Funding That Would Value AI Startup at up to \$340 Billion*, CNBC (Jan. 30, 2025, at 16:15 ET), <https://www.cnbc.com/2025/01/30/openai-in-talks-to-raise-up-to-40-billion-at-340-billion-valuation.html> [<https://perma.cc/VC2C-LMYC>].

²⁰ Hayden Field, *OpenAI Dissolves Team Focused on Long-Term AI Risks, Less Than One Year After Announcing It*, CNBC (May 18, 2024, at 13:49 ET), <https://www.cnbc.com/2024/05/17/openai-superalignment-sutskever-leike.html> [<https://perma.cc/DP8M-BFNN>] (Jan Leike, a cofounder of OpenAI, stated he has “been disagreeing with OpenAI leadership about the company’s core priorities for quite some time, until [they] finally reached a breaking point” and that “[b]uilding smarter-than-human machines is an inherently dangerous endeavor . . . OpenAI is shouldering an enormous responsibility on behalf of all humanity. But over the past years, safety culture and processes have taken a backseat to shiny products.”).

²¹ *Introducing OpenAI*, *supra* note 13 (“Researchers will be strongly encouraged to publish their work, whether as papers, blog posts, or code, and our patents (if any) will be shared with the world. We’ll freely collaborate with others across many institutions and expect to work with companies to research and deploy new technologies.”).

²² Rhiannon Williams & James O’Donnell, *We Finally Have a Definition for Open-Source AI*, MIT TECH. REV. (Aug. 22, 2024), <https://www.technologyreview.com/2024/08/22/1097224/we-finally-have-a-definition-for-open-source-ai/> [<https://perma.cc/D3GD-MMYW>].

its creators, meaning that the model's internal structure and training data is not publicized.²³ Despite starting as an open-source AI model, as OpenAI began to shift its corporate structure and its governance priorities, it has shifted toward a closed-source approach, starting with the development of GPT-3.²⁴ Given the lack of transparency found in the latest closed-source GPT models, it is impossible to know the *exact* content used to train GPT-3 and GPT-4.

B. The Technology Behind OpenAI's Generative AI ChatGPT

Moreover, an understanding of the core claims underlying the *Times*'s lawsuit against OpenAI requires a foundational understanding of what large language models are, how they function, and the source of the datasets used to train them. LLMs are artificial intelligence systems designed to generate humanlike text by recognizing and predicting patterns in language. Understanding this predictive structure is necessary to assess how and why certain data might be used during training.

GPT models are designed to produce natural language text—ranging from sentences and paragraphs to entire documents—in a manner that is coherent and consistent with human language patterns.²⁵ LLMs are “any large-scale language model designed for natural language processing tasks”; hence, GPT models are merely a specific subset of LLMs.²⁶ There is a two-stage training process for GPT models like those created and used by OpenAI: pretraining on extensive datasets, followed by fine-tuning the model for specific tasks.²⁷

²³ Leo Mao, *AI Open and Closed Source: It Is Human Pioneer's Choice*, GWBMA (Apr. 11, 2024), <https://www.registrationchina.com/articles/ai-open-and-closed-source-it-is-human-pioneers-choice/> [<https://perma.cc/Z4GX-GW6C>].

²⁴ Yevgeni Chuvyrov, *Generating Infrastructure-as-Code from Natural Language: Evaluating Fine-Tuned GPT-3 Models for Cloud Infrastructure Provisioning*, STAN. CS224N CUSTOM PROJECT, at 3–6, https://web.stanford.edu/class/archive/cs/cs224n/cs224n.1224/reports/custom_116662212.pdf [<https://perma.cc/D4VL-L58N>] (last visited Aug. 30, 2025).

²⁵ Partha Pratim Ray, *ChatGPT: A Comprehensive Review on Background, Applications, Key Challenges, Bias, Ethics, Limitations and Future Scope*, INTERNET OF THINGS AND CYBER-PHYSICAL SYS., at 121–22 (2023), <https://www.sciencedirect.com/science/article/pii/S266734522300024X> (on file with the Oregon Law Review); *Understanding the Difference Between GPT and LLM: A Comprehensive Comparison*, STACKADEMIC (Oct. 20, 2023) [hereinafter *Difference Between GPT and LLM*], <https://stackademic.com/blog/understanding-the-difference-between-gpt-and-llm-a-comprehensive-comparison-1f624c713507> [<https://perma.cc/RC8P-A7Q4>].

²⁶ *Difference Between GPT and LLM*, *supra* note 25.

²⁷ *Data Labeling: Fine-Tuning vs. Pre-Training: Key Differences for Language Models*, SAPIEN (Oct. 9, 2024), <https://www.sapien.io/blog/fine-tuning-vs-pre-training-key>

During the pretraining phase, GPT models learn to predict the next word in a sentence based on its preceding words or the prompt given by the user.²⁸ This process helps the model learn basic language patterns, including syntax, grammar, and semantics.²⁹ By generalizing these patterns from vast amounts of text, the model develops a relative understanding of human language.³⁰ In other words, because OpenAI's LLMs are trained on an extensive collection of text, the models can readily recognize words and sentences as being organized in sequences with specific dependencies.³¹ This repetitive recognition enables the model to "understand how to cut text into statistical chunks that have some predictability," memorize the patterns within these text blocks, and use this knowledge to predict what the user wants to come next—similar to the autocomplete text feature in email applications.³² LLMs make predictions by using adjustable parameters or weights, which are "complex mathematical transformations that LLMs learn from reading those billions of words," and work to inform the AI "how likely different words or parts of words are to appear together or in a certain order."³³ Once pretraining is complete, the GPT models undergo a fine-tuning process that subjects GPT models "to a more targeted dataset that closely aligns with the specific task at hand," resulting in a more precise output given the user's prompt, purpose, industry, or preference.³⁴

-differences-for-language-models [https://perma.cc/NG39-G57F] ("For example, an LLM can be pre-trained on a massive dataset like Wikipedia to grasp general language patterns and then fine-tuned with customer service scripts to create a chatbot capable of handling customer inquiries with nuanced understanding.").

²⁸ ETHAN MOLLICK, CO-INTELLIGENCE: LIVING AND WORKING WITH AI 2–5 (2024).

²⁹ *See id.* at 5.

³⁰ *Id.* at 5.

³¹ Adam Zewe, *Explained: Generative AI*, MIT NEWS ON CAMPUS & AROUND THE WORLD (Nov. 9, 2023), <https://news.mit.edu/2023/explained-generative-ai-1109> [https://perma.cc/7YRG-22TQ].

³² *Id.*

³³ MOLLICK, *supra* note 28, at 6.

³⁴ Pablo Junco, *The Power of Fine-Tuning in Generative AI*, FORBES (Oct. 10, 2023, at 09:45 ET), <https://www.forbes.com/councils/forbestechcouncil/2023/10/10/the-power-of-fine-tuning-in-generative-ai/> [https://perma.cc/VDJ9-NSUR]; *see also* MOLLICK, *supra* note 28, at 8 ("AI companies hire workers . . . to read AI answers and judge them on various characteristics. In some cases, that might be rating results for accuracy, in others it might be to screen out violent or pornographic answers.").

In 2018, OpenAI introduced its first transformer-based language model, GPT-1.³⁵ The dataset used to train GPT-1 was relatively small compared to the later iterations of GPT models, consisting of only 117 million parameters.³⁶ GPT-1 was trained using the BooksCorpus dataset, a compilation of 7,000 unpublished and self-published books sourced from Smashwords, a platform for such works.³⁷ This dataset provided the language model with exposure to novel data.³⁸

OpenAI released its next iteration, GPT-2, in 2019.³⁹ GPT-2 used a much larger dataset, one consisting of 1.5 billion parameters,⁴⁰ that focused on using web pages curated by humans to both emphasize document quality and produce a larger, more diverse training dataset.⁴¹ Hence, web pages like Reddit, a social media platform, were scraped “to find upvoted articles and [pull] data from all outbound links in the targeted Reddit posts.”⁴²

A year later, OpenAI introduced GPT-3 featuring 175 billion parameters,⁴³ expanding the amount of information in its dataset by using five compilations of data: Common Crawl, WebText2, Books1, Books2, and Wikipedia.⁴⁴ Common Crawl is a dataset that has compiled information from over 250 billion web pages since 2007 and is freely available to researchers.⁴⁵ Another dataset, WebText2, expands on WebText, the dataset that trained GPT-2, and uses a similar

³⁵ Aleksandra Yosifova, *The Evolution of ChatGPT: History and Future*, 365 DATASCIENCE (Aug. 14, 2023), <https://365datascience.com/trending/the-evolution-of-chatgpt-history-and-future/> [https://perma.cc/K5T9-2A9G].

³⁶ Mayankchugh Jobathk, *The Evolution of Large Language Models (LLMs): A Journey from GPT to GPT-40*, MEDIUM (July 18, 2024), <https://medium.com/@mayankchugh.jobathk/the-evolution-of-large-language-models-llms-a-journey-from-gpt-to-gpt-40-618765889c98> [https://perma.cc/XES7-7PAN].

³⁷ Art Neill, James Thomas & Erika Lee, *A Framework for Applying Copyright Law to the Training of Textual Generative Artificial Intelligence*, 32 TEX. INTELL. PROP. L.J. 225, 231 (2024).

³⁸ Priya Shree, *The Journey of Open AI GPT Models*, MEDIUM (Nov. 9, 2020), <https://medium.com/walmartglobaltech/the-journey-of-open-ai-gpt-models-32d95b7b7fb2> [https://perma.cc/7QGH-REGW].

³⁹ Jobathk, *supra* note 36.

⁴⁰ *Id.*

⁴¹ Alec Radford et al., *Language Models Are Unsupervised Multitask Learners*, OPENAI at 3, https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf [https://perma.cc/42RD-UGP4] (last visited Aug. 20, 2025).

⁴² Neill, Thomas & Lee, *supra* note 37, at 231.

⁴³ Jobathk, *supra* note 36.

⁴⁴ Neill, Thomas & Lee, *supra* note 37, at 231–32.

⁴⁵ See *Frequently Asked Questions*, COMMON CRAWL, <https://commoncrawl.org/faq> [https://perma.cc/2X9F-7M7H] (last visited Mar. 2, 2025).

selection criterion, targeting websites linked to Reddit posts with at least three upvotes.⁴⁶ In contrast, the Books1 and Books2 datasets are less precisely defined but appear to consist of books in the public domain.⁴⁷ The final dataset, Wikipedia, seemingly includes all text and data available on the popular and free-content encyclopedia platform.⁴⁸

Finally, OpenAI released its most recent version, GPT-4, in 2023, with a staggering estimation of 1 trillion parameters.⁴⁹ Hence, GPT-4 contains the most expansive training dataset yet, containing “web texts, books, news articles, social media posts, code snippets, and other unspecified sources.”⁵⁰ However, unlike previous iterations of GPT, OpenAI has become a closed-source system, meaning it has completely “closed off” information about the types of data used for the training of this version of GPT.⁵¹

Despite the closed-source system of OpenAI’s latest GPT models, a typical dataset is comprised of text sourced from “the internet, public domain books and research articles, and assorted other free sources of content that researchers can find.”⁵² However, as AI development continues to advance, the demand for high-quality training data has become a pressing issue.⁵³ Many AI companies are rapidly depleting readily accessible, high-value sources, prompting concerns over the increasing likelihood that training datasets contain copyrighted material, whether inadvertently or deliberately.⁵⁴

C. The Copyrightability of GPT-4’s Training Dataset

Furthermore, distinguishing between public domain works, openly licensed materials, and fully copyrighted texts is critical to evaluate

⁴⁶ Roger Montti, *How to Block OpenAI ChatGPT from Using Your Website Content*, SEARCH ENGINE J. (Feb. 2, 2023), <https://www.searchenginejournal.com/how-to-block-chatgpt-from-using-your-website-content/478384/> [<https://perma.cc/K2UM-46RK>].

⁴⁷ Gregory Roberts, *AI Training Datasets: The Books1+Books2 That Big AI Eats for Breakfast*, GREGOREITE (Dec. 14, 2022), <https://gregoreite.com/drilling-down-details-on-the-ai-training-datasets/> [<https://perma.cc/28F2-9QAR>].

⁴⁸ *Id.*

⁴⁹ Jobathk, *supra* note 36.

⁵⁰ Neill, Thomas & Lee, *supra* note 37, at 232.

⁵¹ *Id.*

⁵² MOLLICK, *supra* note 28, at 7.

⁵³ *Id.* at 7–8.

⁵⁴ *Id.* at 7.

whether the use of such content by models like GPT-4 (the subject of this lawsuit) aligns with or violates existing copyright law.⁵⁵

A significant portion of the training data for LLMs is likely derived from public domain works. Works in the public domain are no longer protected by copyright law and instead are owned by the public, subject to the public's use without permission.⁵⁶ Although the precise nature of information in the datasets Books1 and Books2 is undisclosed, these datasets are widely speculated to consist of, at least in part, books that have entered the public domain.⁵⁷

Wikipedia, on the other hand, is published under an open license, allowing users to access Wikipedia's articles for free and for reuse.⁵⁸ Aside from specific quoted content that may fall under separate copyright protections, Wikipedia is generally licensed under the Creative Commons Attribution-ShareAlike 4.0 International License and the GNU Free Documentation License.⁵⁹ As a result, Wikipedia's content may typically be used without risk of copyright infringement provided that proper attribution is given, consistent with the terms of these open licenses.⁶⁰

As for the last category of sources, fully protected copyrighted works not subject to open licensing, it is more unclear as to the exact amount of copyrighted material that went into the training of LLMs like GPT-4. At the very least, we know that some copyrighted materials went into the training of GPT-1 because BooksCorpus, one of the

⁵⁵ *Id.* at 233.

⁵⁶ Richard Stim, *Welcome to the Public Domain*, COPYRIGHT & FAIR USE STANFORD LIBRS., <https://fairuse.stanford.edu/overview/public-domain/welcome/#:~:text=The%20term%20%E2%80%9Cpublic%20domain%E2%80%9D%20refers,one%20can%20ever%20own%20it> [https://perma.cc/4VPL-2LT9] (last visited Sep. 17, 2025) (Public domain “refers to creative materials that are not protected by intellectual property laws such as copyright, trademark, or patent laws. The public owns these works, not an individual author or artist. Anyone can use a public domain work without obtaining permission, but no one can ever own it.”).

⁵⁷ Roberts, *supra* note 47; see also Jack Bandy & Nicholas Vincent, *Addressing “Documentation Debt” in Machine Learning Research: A Retrospective Datasheet for BookCorpus*, <https://arxiv.org/pdf/2105.05241> [https://perma.cc/82DB-EBE6] (May 11, 2021).

⁵⁸ Yana Welinder, *Free as in Open Access and Wikipedia*, THE CTR. FOR INTERNET & SOC'Y (Oct. 20, 2014, at 03:00 PT), <https://cyberlaw.stanford.edu/blog/2014/10/free-open-access-and-wikipedia/> [https://perma.cc/BS2B-PTCE].

⁵⁹ *Wikipedia: FAQ/Copyright*, WIKIPEDIA, https://en.wikipedia.org/wiki/Wikipedia:FAQ/Copyright#Can_I_reuse_Wikipedia's_content_somewhere_else [https://perma.cc/9PSX-NKQY] (last visited Sep. 17, 2025).

⁶⁰ *Id.*

sources in the dataset, also contained copyrighted books.⁶¹ Moreover, authors of the copyrighted works in BooksCorpus were unable “to opt out of the inclusion of their works” into the GPT training datasets.⁶² Therefore, because existing datasets contained works that are fully subject to copyright law restrictions, the *exact* number of such copyrighted works in the dataset is irrelevant; the fact that at least one copyrighted work made it into the dataset is sufficient.

II

THE *TIMES*’S ALLEGATIONS

In late December 2023, *The New York Times* filed an action against defendants OpenAI, Inc. and Microsoft Corporation alleging seven counts of unlawful use and infringement of the *Times*’s articles, with its lead claim being direct copyright infringement under 17 U.S.C. § 501.⁶³

The central claim⁶⁴ is that OpenAI has infringed upon the *Times*’s copyrights through the unlicensed and unauthorized use and reproduction of the *Times*’s content in the training of its GPT models.⁶⁵ First, the complaint alleges that LLMs like GPT are capable of “memorizing” or “regurgitating” portions of works included in their training data.⁶⁶ When this occurs, the models can generate near-verbatim reproductions of the original works the models were trained on, such as the *Times*’s copyrighted articles in the dataset.⁶⁷ Second, and relatedly, the LLMs can produce synthetic search results that, when prompted, reproduce “significantly more expressive content from [an]

⁶¹ See *supra* text accompanying notes 50–54.

⁶² Neill, Thomas & Lee, *supra* note 37, at 234; see also Bandy & Vincent, *supra* note 57.

⁶³ New York Times Complaint, *supra* note 4, at 60–67.

⁶⁴ The *Times* Complaint asserts seven counts against Microsoft Corporation and OpenAI, Inc. and its related business entities. The following Counts are alleged against all defendants: Copyright Infringement under 17 U.S.C. § 501, Contributory Copyright Infringement, Violation of the Digital Millennium Copyright Act under 17 U.S.C. § 1202, Common Law Unfair Competition by Misappropriation, and Trademark Dilution under 15 U.S.C. § 1125(c). The *Times* also alleges Vicarious Copyright Infringement against Microsoft, OpenAI, Inc., OpenAI LP, OAI Corporation LLC, OpenAI Holdings LLC, and OpenAI Global LLC. The *Times* also has a separate Contributory Copyright Infringement claim against just Microsoft Corporation. For the purposes of this Comment, I focus on only Count I: Copyright Infringement under § 501 as it pertains to defendant OpenAI, Inc.

⁶⁵ New York Times Complaint, *supra* note 4, at 2.

⁶⁶ *Id.* at 23–24.

⁶⁷ *Id.*

original article than what would traditionally be displayed” through an online search engine.⁶⁸ As the *Times* alleges, this ability enables users to circumvent the *Times*’s paywall and undermine its subscription-based business model.⁶⁹ At its core, the *Times* stresses that these practices gravely threaten high-quality journalism,⁷⁰ because if OpenAI users can easily generate summaries or near-verbatim reproductions of the *Times*’s articles, at no cost or consequence, it would “obviate the need” to purchase access to the *Times* directly, thereby undermining its business model.⁷¹

The complaint also emphasizes that, given the need for high-quality textual works in GPT datasets,⁷² the algorithm was designed to disproportionately curate the *Times*’s works.⁷³ The complaint alleges that the WebText dataset was created with the intention of securing high-quality documents and, thus, contains “a staggering amount of scraped content from the *Times*,” with the *Times*’s domain name being “one of the ‘top 15 domains by volume’ in the WebText dataset.”⁷⁴ As the complaint alleges, because OpenAI acknowledges that “datasets [it] view[s] as higher-quality are sampled more frequently” during the training process, by its own admission, it suggests that high-quality content, such as material from the *Times*, played a more significant and valuable role in training GPT models compared to lower-quality sources.⁷⁵

Ultimately, the complaint underscores the societal stakes, stating that “[i]ndependent journalism is vital to our democracy” and is “increasingly rare and valuable.”⁷⁶ In contrast, the *Times* characterizes OpenAI as merely a “multi-billion-dollar for-profit business built in large part on the unlicensed exploitation of copyrighted works” that shed its nonprofit status a mere three years after its founding.⁷⁷ The *Times*’s complaint suggests that OpenAI obtained public trust through

⁶⁸ *Id.* at 40, 42, 44, 46.

⁶⁹ *Id.* at 59.

⁷⁰ *Id.* at 14.

⁷¹ *Id.* at 22.

⁷² See *supra* text accompanying notes 50–54.

⁷³ New York Times Complaint, *supra* note 4, at 20.

⁷⁴ *Id.* at 25.

⁷⁵ *Id.* at 27.

⁷⁶ *Id.* at 1.

⁷⁷ *Id.* at 17.

its brief nonprofit and open-source phase to, within a few years, become what it once promised never to be.⁷⁸

III OPENAI'S DEFENSES

Defendant OpenAI argues that the crux of the *Times*'s infringement allegations rest on "training data regurgitation," which are merely "uncommon and unintended phenomena" that OpenAI has already identified as a problem and is actively working to resolve.⁷⁹ OpenAI even goes as far as accusing the *Times* of "pa[ying] someone to hack" its products, claims the *Times* engaged in "deceptive prompt[ing]" to exploit a "bug," and states the *Times* was able to obtain the verbatim outputs in Exhibit J of the Complaint only by engaging in "tens of thousands of attempts."⁸⁰ This, OpenAI argues, violates its terms of use and is not how "[n]ormal people" use its products.⁸¹

More importantly, OpenAI argues that its use of the *Times*'s articles for the training of its GPT models constitutes fair use.⁸² In particular, OpenAI argues that it uses the copyrighted content as a mere cog in the technological process and development of its LLMs to produce "new, different, and innovative products."⁸³

Moreover, OpenAI seeks to partially dismiss the direct infringement claim, arguing the claim is time-barred by the three-year limitations period because OpenAI's use of the *Times*'s articles to train the LLMs occurred more than three years ago.⁸⁴ However, as the *Times* notes in its opposition brief, "OpenAI's statute-of-limitations argument is narrow, addressing only the training of GPT-2 in 2019 and GPT-3 in 2020 [i]t does not apply to the 'orders of magnitude more

⁷⁸ See *Introducing OpenAI*, *supra* note 13; Chloe Xiang, *OpenAI Is Now Everything It Promised Not to Be: Corporate, Closed-Source, and For-Profit*, VICE (Feb. 28, 2023, at 12:35 PT), <https://www.vice.com/en/article/openai-is-now-everything-it-promised-not-to-be-corporate-closed-source-and-for-profit/> [<https://perma.cc/S7P8-JSTJ>].

⁷⁹ Defendant's Memorandum in Support of Motion to Dismiss, *supra* note 9, at 10–11.

⁸⁰ *Id.* at 2. Excerpts of Exhibit J in Plaintiff's complaint show a side-by-side comparison of ChatGPT-4's outputs and various *New York Times* articles; *see also id.* at 30–37.

⁸¹ *Id.*

⁸² *OpenAI and Journalism*, OPENAI (Jan. 8, 2024), <https://openai.com/index/openai-and-journalism/> (on file with the Oregon Law Review); *see also* Bobby Allyn, *The New York Times' Takes OpenAI to Court. ChatGPT's Future Could Be on the Line*, NPR (Jan. 14, 2025, at 16:27 ET), <https://www.npr.org/2025/01/14/nx-s1-5258952/new-york-times-openai-microsoft> (on file with the Oregon Law Review).

⁸³ Defendant's Memorandum in Support of Motion to Dismiss, *supra* note 9, at 3.

⁸⁴ *Id.* at 15.

powerful’ GPT-3.5 and GPT-4 models developed in 2022 and 2023.”⁸⁵ Interestingly, prior to filing a formal legal response, OpenAI publicly responded to the *Times*’s complaint in a blog post delineating four main points: (1) OpenAI “collaborate[s] with news organizations,” (2) OpenAI’s use of copyrighted materials to train its LLMs constitutes fair use and is hence not infringement, (3) regurgitation is a rare phenomenon that it is “working to drive to zero,” and (4) the regurgitations the *Times* argues as a basis for its infringement claim are so out of the ordinary that it raises inferences that the *Times* “either instructed the model to regurgitate or cherry-picked their examples from many attempts.”⁸⁶

IV

MERITS OF THE LEGAL CLAIMS

While the legal dispute is likely to turn on the court’s application of the fair use defense, the *Times* will still have to establish direct copyright infringement. Where it is successful in demonstrating direct infringement, it will similarly need to show that OpenAI’s fair use defense is inapplicable in this case. Given the relatively low bar for establishing ownership and unauthorized copying,⁸⁷ and the available evidence suggesting the *Times*’s content was used in training datasets, the *Times* is likely to succeed in making out a prima facie case of copyright infringement.⁸⁸ As a result, the disposition of the case will likely hinge on OpenAI’s fair use defense argument. This Section IV.B offers a detailed application of the four statutory fair use factors.

A. Direct Copyright Infringement

Under the Copyright Act of 1976, 17 U.S.C. § 101, the owner of a copyrighted work has the exclusive right to reproduce, perform publicly, display publicly, prepare derivative works of, and distribute copies of their copyrighted work.⁸⁹ Where a defendant violates a plaintiff’s exclusive rights in their copyrighted work, the defendant is

⁸⁵ Plaintiff’s Memorandum in Opposition to OpenAI Defendants’ Partial Motion to Dismiss at 6, *N.Y. Times Co. v. Microsoft Corp.*, No. 23-CV-11195 (S.D.N.Y. Apr. 4, 2025).

⁸⁶ *OpenAI and Journalism*, *supra* note 82.

⁸⁷ *See, e.g.*, 17 U.S.C. § 410(c) (“[T]he certificate of a registration made before or within five years after first publication of the work shall constitute prima facie evidence of the validity of the copyright . . .”).

⁸⁸ Defendant’s Memorandum in Support of Motion to Dismiss, *supra* note 9, at 7.

⁸⁹ 17 U.S.C. § 106.

liable for damages.⁹⁰ Therefore, for the *Times* to prevail on its direct copyright infringement claim against OpenAI, it must establish by a preponderance of the evidence that (1) it owns a valid copyright, and (2) the defendant copied “constituent elements of the work that are original.”⁹¹ The second element requires showing both that OpenAI actually copied the work and that its copy was substantially similar to the work.⁹²

Copyright registrations are “prima facie evidence of the validity of the copyright” if “made before or within five years after first publication of the work.”⁹³ Given that the *Times* has registered its articles, it is unlikely OpenAI will be successful in disputing this element.⁹⁴

Hence, the major legal disputes are likely to occur under the second element. To satisfy the second element of a copyright claim, the *Times* must demonstrate that OpenAI copied its work and that such copying was unlawful because OpenAI’s allegedly infringing work is substantially similar to the protectable elements of the *Times*’s copyrighted works.⁹⁵ Actual copying means that “the defendant did, in fact, use the copyrighted work in creating his own.”⁹⁶ Because direct evidence of copying is difficult to establish, the *Times* can rely on circumstantial evidence of copying “‘by demonstrating that the person who composed the defendant’s work had access to the copyrighted material,’ and that there are similarities between the two works that are ‘probative of copying.’”⁹⁷ The *Times* will likely be able to establish that OpenAI had access to its copyrighted articles because, in July of 2020, “OpenAI had disclosed that the *Times*’s articles were a tiny part of the diverse datasets that had been used to train these language models[,]” establishing that OpenAI undoubtedly had access to such

⁹⁰ Warren v. John Wiley & Sons, Inc., 952 F. Supp. 2d 610, 616 (S.D.N.Y. 2013).

⁹¹ Disney Enter., Inc. v. Sarelli, 322 F. Supp. 3d 413, 441 (S.D.N.Y. 2018) (quoting Muller v. Twentieth Century Fox Film Corp., 794 F. Supp. 2d 429, 439 (S.D.N.Y. 2011)).

⁹² Dam Things from Den. v. Russ Berrie & Co., Inc., 290 F.3d 548, 561–62 (3d Cir. 2002).

⁹³ 17 U.S.C. § 410(c).

⁹⁴ Copyright Notice, THE NEW YORK TIMES, <https://help.nytimes.com/hc/en-us/articles/115014792127-Copyright-Notice> [<https://perma.cc/6JF8-CPC2>] (last visited Sept. 13, 2024).

⁹⁵ Sarelli, 322 F. Supp. 3d at 442.

⁹⁶ Tanksley v. Daniels, 902 F.3d 165, 173 (3d Cir. 2018).

⁹⁷ Jorgensen v. Epic/Sony Recs., 351 F.3d 46, 51 (2d Cir. 2003) (first quoting Herzog v. Castle Rock Ent., 193 F.3d 1241, 1249 (11th Cir. 1999); and then quoting Repp v. Weber, 132 F.3d 882, 889 (2d Cir. 1997)).

copyrighted works.⁹⁸ Additionally, because the *Times* is popular worldwide, it is unlikely that “access” will be a major issue in this case.

Moreover, the *Times* must still prove that the similarities between its copyrighted articles and OpenAI’s outputs are so similar that an “average lay observer would recognize the alleged copy as having been appropriated from the copyrighted work” under the substantial similarity analysis.⁹⁹ Because the case here involves works of literary expression, the *Times* must demonstrate that the similarities shared by its articles and OpenAI’s outputs are “something more than mere generalized ideas or themes.”¹⁰⁰ It is likely that the *Times* will be able to establish a substantial similarity to OpenAI’s outputs—at least in terms of Exhibit J in the complaint, which introduced one hundred examples¹⁰¹ where ChatGPT was able to nearly regurgitate parts of the *Times*’s articles verbatim and provide summaries that are “significantly longer and more detailed” than what is traditionally accessible through normal search engines.¹⁰²

B. Application of the Fair Use Defense

The fair use defense operates as a limitation on a copyright owner’s rights of exclusivity to promote the breathing room necessary to promote creative expression.¹⁰³ The fair use defense to copyright infringement is a question of law that is determined by underlying facts on a case-by-case basis.¹⁰⁴ Hence, although it is not possible to categorically assert that every aspect of the LLMs’ training are per se protected under a fair use defense, the discussion below examines how the training of generative AI aligns with existing fair use jurisprudence.

Under 17 U.S.C. § 107, the fair use defense requires consideration of

- (1) the purpose and character of the use, including whether such use is of a commercial nature or is for nonprofit educational purposes;

⁹⁸ Defendant’s Memorandum in Support of Motion to Dismiss, *supra* note 9, at 7 (OpenAI admits the *Times*’s articles exist in the training dataset but argues their use is protected under the fair use doctrine).

⁹⁹ *E.g.*, *Sarelli*, 322 F. Supp. 3d at 442 (quoting *Blakeman v. The Walt Disney Co.*, 613 F. Supp. 2d 288, 304 (E.D.N.Y. 2009); *see also* *Warner Bros. Inc. v. Am. Broad. Cos., Inc.*, 654 F.2d 204, 208 (2d Cir. 1981) (“In the case of literary works, it is axiomatic that copyright protection only extends to the expression of the author’s idea, not to the idea itself.”).

¹⁰⁰ *Warner Bros. Inc.*, 654 F.2d at 208.

¹⁰¹ New York Times Complaint, *supra* note 4, at 2.

¹⁰² *Id.* at 3.

¹⁰³ *Green v. U.S. Dep’t of Just.*, 111 F.4th 81, 101 (D.C. Cir. 2024).

¹⁰⁴ *See* *Google LLC v. Oracle Am., Inc.*, 593 U.S. 1, 23 (2021).

(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 value . . . of the copyrighted work.¹⁰⁵

I. Purpose and Character

Although courts must consider all four fair use factors, the “purpose and character” of the use has historically played a central role in the analysis.¹⁰⁶ This factor addresses “[t]he heart of the fair use inquiry”—namely, whether the defendant’s use is “transformative,” and whether it is for commercial or nonprofit educational purposes.¹⁰⁷ A use is transformative if it uses the original work to create something new, with a different purpose or character, rather than merely superseding the original.¹⁰⁸ Conversely, a use is not transformative where it does not serve a “further purpose or different character” than the original work, meaning that the use merely acts as a substitute or replacement for the original work.¹⁰⁹

17 U.S.C. § 107 provides an inexhaustive list of potentially fair-use, transformative purposes—“criticism, comment, news reporting, teaching . . . scholarship, or research”¹¹⁰—all of which indicate uses that are incapable of supplanting the original works’ use.¹¹¹ However, given that most secondary uses add something at least slightly novel, this factor asks “‘whether and to what extent’ the use at issue has a

¹⁰⁵ 17 U.S.C. § 107.

¹⁰⁶ Nat’l Acad. of Television Arts and Scis., Inc. v. Multimedia Sys. Design, Inc., 551 F. Supp. 3d 408, 421–22 (S.D.N.Y. 2021).

¹⁰⁷ *Id.*

¹⁰⁸ Richard Stim, *Measuring Fair Use: The Four Factors*, COPYRIGHT & FAIR USE STANFORD LIBR., <https://fairuse.stanford.edu/overview/fair-use/four-factors/> [<https://perma.cc/84P2-BXGW>] (last visited Sept. 14, 2025).

¹⁰⁹ *Andy Warhol Found. for the Visual Arts, Inc. v. Goldsmith*, 598 U.S. 508, 525 (2023).

¹¹⁰ 17 U.S.C. § 107.

¹¹¹ *Andy Warhol*, 598 U.S. at 528 (“Criticism of a work, for instance, ordinarily does not supersede the objects of, or supplant, the work. Rather, it uses the work to serve a distinct end.”); *see also* *Folsom v. Marsh*, 9 F. Cas. 342, 344–45 (C.C.D. Mass. 1841) (No. 4,901,2) (“[I]t is as clear, that if he thus cites the most important parts of the work, with a view, not to criticise, but to supersede the use of the original work, and substitute the review for it, such a use will be deemed in law a piracy.”). *But see* *Harper & Row Publishers, Inc. v. Nation Enters.*, 471 U.S. 539, 561 (1985) (“The fact that an article arguably is ‘news’ and therefore a productive use is simply one factor in a fair use analysis.”).

purpose or character different from the original.”¹¹² In sum, the first fair use factor weighs the degree to which the secondary use of a copyrighted work serves a distinct purpose or transforms the original in some meaningful way against the commercial nature of the use.¹¹³ When a secondary use closely mirrors the purpose of the original work and is driven by profit, the first factor will generally weigh against a finding of fair use, unless compelling justification for the copying exists.¹¹⁴

Importantly, a secondary use that merely adds new expression, meaning, or message—but retains the same purpose as the original—does not qualify as transformative.¹¹⁵ In *Andy Warhol*, a magazine hired Goldsmith to photograph the musician Prince.¹¹⁶ Years later, Andy Warhol used one of Goldsmith’s copyrighted photographs of Prince without permission as the basis for a series of colorful prints, one of which was licensed and published on the cover of a popular magazine.¹¹⁷ Despite their stylistic differences, the Court found that Warhol’s prints served the same commercial purpose as Goldsmith’s original photographs: to illustrate articles about Prince in magazines.¹¹⁸ Therefore, although the two works “were not perfect substitutes” of each other, because both works competed in the same licensing market, the Court held that Warhol’s use was not transformative and weighed against a finding of fair use.¹¹⁹

Conversely, where a use does not replicate the original copyrighted work but merely furnishes a searchable database and provides information about the original work, the use is likely transformative.¹²⁰ In *Authors Guild*, Google scanned millions of books for its Google Books project, including works in the public domain and copyrighted materials, without the consent of the owners of such materials.¹²¹ Google kept the original scans of these books, but allowed its users to

¹¹² *Andy Warhol*, 598 U.S. at 529 (quoting *Campbell v. Acuff-Rose Music, Inc.*, 510 U.S. 569, 579 (1994)) (“The larger the difference, the more likely the first factor weighs in favor of fair use.”).

¹¹³ *Id.*; see also *Google LLC v. Oracle Am., Inc.*, 593 U.S. 1, 29 (2021).

¹¹⁴ *Andy Warhol*, 598 U.S. at 531–32.

¹¹⁵ *Id.* at 550.

¹¹⁶ *Id.* at 516.

¹¹⁷ *Id.* at 518–19.

¹¹⁸ *Id.* at 550.

¹¹⁹ *Id.* at 536–37.

¹²⁰ *Authors Guild v. Google, Inc.*, 804 F.3d 202, 219 (2d Cir. 2015).

¹²¹ *Id.* at 208.

access only a limited view of the text in response to queries.¹²² The plaintiffs challenged Google on the basis of two uses: (1) its use in providing a search engine for users to determine whether a key term appeared in a book, and if so, how many times, and (2) its use of the copyrighted works to allow users to see a “snippet” or a portion of the original text.¹²³ Because Google copied the original books for the purpose of making information about those books available to its users, rather than serving merely as a substitute for the books, the court held that Google’s use was transformative because it fundamentally altered the purpose of the original works.¹²⁴ Moreover, this transformation conferred significant public benefits, such as facilitating scholarship and expanding access to knowledge, without supplanting the market for the original books.¹²⁵ Hence, although Google derived a commercial benefit from the project through advertising, the court found that this factor did not outweigh the transformative nature of the use.¹²⁶

Similar to *Andy Warhol*, and distinct from *Authors Guild*, OpenAI’s use of the *Times*’s articles to train its generative AI tools does not alter the fundamental purpose of the original works. The *Times*’s purpose in the creation of its articles centers around investigative reporting, breaking news reporting, beat reporting, reviews and analysis, and commentary and opinion.¹²⁷ Arguably, OpenAI’s use is to train its LLMs for the ultimate purpose of providing reports, analysis, commentary, and opinion-based outputs to its users.¹²⁸ Therefore, just as Warhol’s prints were not “perfect substitutes” for Goldsmith’s photograph¹²⁹ but nevertheless shared a similar purpose and commercial function, OpenAI’s ingestion of the *Times*’s articles serves

¹²² *Id.* at 209.

¹²³ *Id.* at 222.

¹²⁴ *Id.* at 218.

¹²⁵ *See id.* at 207.

¹²⁶ *Id.* at 219.

¹²⁷ New York Times Complaint, *supra* note 4, at 10–12.

¹²⁸ *See How ChatGPT and Our Foundation Models Are Developed*, OPENAI, <https://help.openai.com/en/articles/7842364-how-chatgpt-and-our-foundation-models-are-developed> (on file with the Oregon Law Review) (last visited June 27, 2025) (Users can “use ChatGPT for a wide range of tasks, including organizing and summarizing information, assisting with translations, analyzing or generating images, inspiring creativity and ideas, and other everyday activities.”).

¹²⁹ *Andy Warhol Found. for the Visual Arts, Inc. v. Goldsmith*, 598 U.S. 508, 536 (2023).

a purpose closely aligned with that of the source material: the communication of information.

Moreover, in contrast to *Authors Guild*, OpenAI's use fundamentally differs in terms of function and output. First, whereas Google Books was a tool for locating information within a book, thus not acting as an actual substitute for the information itself,¹³⁰ OpenAI's LLMs ingest the expressive content of the *Times*'s articles for the purpose of generating new textual outputs that resemble, and at times actually copy the original work's substance, style, and overall expression.¹³¹ Second, unlike *Authors Guild*, GPT's outputs go beyond what is revealed in typical snippets or ordinary search results.¹³² Even when these AI-generated summaries include links or citations to the original source, users have fewer incentives to actually navigate to those sources because the expressive content, either quoted or paraphrased, is already incorporated into the narrative response.¹³³ In fact, the inclusion of attribution may inadvertently increase user trust in the summary itself, reducing the likelihood that users will click through to the original article.¹³⁴ As a result, users may never visit the *Times*'s website at all, directly undermining its ability to monetize access to its content—an ability the Copyright Clause is intended to provide.¹³⁵

Furthermore, OpenAI's status as a for-profit enterprise undermines its fair use defense. The commercial nature of OpenAI's activities, especially its monetization of the GPT product line through premium subscriptions, further undermines its fair use defense.¹³⁶ Therefore, given the similarity and near synonymous purpose of the parties' use, in combination with the commercial nature of the use, it is likely that the purpose and character factor will weigh against the application of the fair use defense.

On the other hand, OpenAI may argue that its purpose is narrower, using the *Times*'s articles to merely train its LLMs to understand

¹³⁰ *Authors Guild v. Google, Inc.*, 804 F.3d 202, 224 (2d Cir. 2015).

¹³¹ See New York Times Complaint, *supra* note 4, at 30–37.

¹³² See *id.* at 2.

¹³³ *Id.* at 37.

¹³⁴ *Id.*

¹³⁵ U.S. CONST. art. I, § 8, cl. 8; see also *Harper & Row Publishers, Inc. v. Nation Enters.*, 471 U.S. 539, 562 (1985) (“[E]very commercial use of copyrighted material is presumptively an unfair exploitation of the monopoly privilege that belongs to the owner of the copyright.”) (internal quotations omitted).

¹³⁶ See *Introducing ChatGPT Plus*, OPENAI (Feb. 1, 2023), <https://openai.com/index/chatgpt-plus/> (on file with the Oregon Law Review).

the mechanics and basic rules of human language rather than to inform readers about current events.¹³⁷ While this distinction may be technically correct at the training stage, the ultimate purpose of the use is to generate outputs readable by humans that can and have regurgitated content in its dataset verbatim, particularly that of the *Times*'s content.¹³⁸ Under such circumstances, the purported distinction between “training” and “informing” collapses. The former is instrumental to the latter, and OpenAI's models, trained on the *Times*'s content, fulfill the same market facing role as the original works. Thus, the purpose and character of OpenAI's uses are insufficiently transformative and weigh against a finding of fair use.

2. Nature of the Copyrighted Works

The second factor, the nature of the copyrighted work, evaluates where the original work lies on the spectrum between factual and creative, with the underlying premise that creative works “are closer to the core of intended copyright protection.”¹³⁹ The *Times*'s articles, as news reporting, generally reside closer to the factual end of the spectrum that, again, can favor fair use. However, copyright law nevertheless safeguards the author's particular *expression* of those facts, including the selection of facts, the ordering of facts, and the overall arranging of facts.¹⁴⁰

Facts as individual propositions exist within the public domain; however, when an author compiles those facts in a way that reflects their original creativity, judgment, and effort, the protection of copyright law is more likely to become available.¹⁴¹ In *Schroeder*, the plaintiff compiled a book listing the names and addresses of suppliers of gardening items and tools.¹⁴² The defendant copied the names and

¹³⁷ Defendant's Memorandum in Support of Motion to Dismiss, *supra* note 9, at 4–5.

¹³⁸ See New York Times Complaint, *supra* note 4, at 30–37.

¹³⁹ *Campbell v. Acuff-Rose Music, Inc.*, 510 U.S. 569, 586 (1994); see also *Hachette Book Grp., Inc. v. Internet Archive*, 115 F.4th 163, 187 (2d Cir. 2024).

¹⁴⁰ *Authors Guild v. Google, Inc.*, 804 F.3d 202, 220 (2d Cir. 2015) (“The mere fact that the original is a factual work therefore should not imply that others may freely copy it. Those who report the news undoubtedly create factual works. It cannot seriously be argued that, for that reason, others may freely copy and re-disseminate news reports.”); *Harper & Row*, 471 U.S. at 561 (“The fact that an article arguably is ‘news’ and therefore a productive use is simply one factor in a fair use analysis.”).

¹⁴¹ *Schroeder v. William Morrow & Co.*, 566 F.2d 3, 5–6 (7th Cir. 1977) (“Another is entitled to make his own compilation of the same names and addresses, using information in the public domain, but he is not entitled merely to copy the copyrighted list.”).

¹⁴² *Id.* at 4.

addresses from almost fifty percent of the plaintiff's book.¹⁴³ Even though the defendant copied only the names and addresses from the plaintiff's book and not any of "the accompanying descriptive material," the court rejected the fair use defense on the ground that copyright protection extends to "'the selection, the ordering and arrangement' of the names and addresses."¹⁴⁴

Here, it is unlikely OpenAI's use does not merely involve the extraction of bare facts or general ideas available in the public domain.¹⁴⁵ Instead, the *Times* presents evidence that OpenAI's models can reproduce entire paragraphs from its articles, capturing not just factual content, but the original phrasing, sequencing, and stylistic elements that reflect significant editorial judgment.¹⁴⁶ If the mere compilation of names and addresses is sufficient to reject the fair use defense, surely the *Times*'s worldwide, three-time Pulitzer-award-winning articles are entitled to copyright protection,¹⁴⁷ even if the main components are facts.

Furthermore, OpenAI's use does not merely involve isolated factual information or general ideas drawn from the *Times*'s reporting. Instead, OpenAI's models, as shown in Exhibit J,¹⁴⁸ can reproduce substantial portions of the *Times*'s articles, including entire paragraphs.¹⁴⁹ The copying of lengthy, coherent portions of original text evidence a direct taking of the author's expressive choices, regardless of the work's fictional or nonfictional status.

Moreover, the court will consider whether the original work is published, in which an unpublished original work is more likely to weigh against a finding of fair use.¹⁵⁰ The *Times*'s articles are published works; however, the mere fact of publication alone is insufficient to satisfy the second factor in the face of a defendant's wholesale copying of entire paragraphs of narrative text.

¹⁴³ *Id.* at 4–5.

¹⁴⁴ *Id.* at 6.

¹⁴⁵ See *supra* text accompanying notes 50–54.

¹⁴⁶ See New York Times Complaint, *supra* note 4, at 30–37.

¹⁴⁷ See *Awards and Recognition*, N.Y. TIMES, <https://www.nytc.com/> [<https://perma.cc/JR4T-83G5>].

¹⁴⁸ See New York Times Complaint, *supra* note 4, at 30–37.

¹⁴⁹ Stim, *supra* note 56.

¹⁵⁰ *Authors Guild v. Google, Inc.*, 804 F.3d 202, 220 (2d Cir. 2015).

Ultimately, copyright law does not require absolute novelty or inventiveness—only a minimal degree of creativity.¹⁵¹ Thus, the *Times* need not demonstrate that its reporting introduces new facts or revolutionary ideas; rather, it must show that its articles embody original expressive choices in how facts are selected, framed, and communicated to the public. Therefore, even if the factual and published nature somewhat reduces the *Times*'s argument, the overall fair use argument will likely still weigh in the *Times*'s favor because the purpose of OpenAI's uses is nearly synonymous with the *Times*'s, and therefore it can act as a direct market substitute.

3. Amount and Substantiality

The third factor asks whether “‘the amount and substantiality of the portion used in relation to the copyrighted work as a whole’ are reasonable in relation to the copying’s purpose.”¹⁵² Under this factor, where the most important parts of the original work are copied or are largely copied, a court is less likely to find that this factor weighs in favor of fair use because it increases the likelihood that the defendant’s work is an effective substitute for the original.¹⁵³

Even where the infringement is quantitatively small, if the infringement captures the most valuable and expressive portion of the original copyrighted work—in other words, the “heart” of the work—the infringement is considered substantial under the third factor.¹⁵⁴ In *Harper & Row*, former President Ford contracted with a magazine the exclusive right to publish his memoirs that contained his reflections on the Watergate crisis and his pardon of former President Nixon.¹⁵⁵ Before the magazine could publish excerpts of the memoir, a political commentary magazine, *The Nation*, obtained an unauthorized copy of the memoir and published roughly 300 words of copyrighted material.¹⁵⁶ As a result of *The Nation*'s article, the contracted magazine

¹⁵¹ *Feist Publ’ns, Inc. v. Rural Tel. Serv. Co., Inc.*, 499 U.S. 340, 358–59 (1991) (holding copyright protects original works that possess a minimal level of creativity regardless of whether the underlying information is fact or fiction).

¹⁵² *Campbell v. Acuff-Rose Music, Inc.*, 510 U.S. 569, 570 (1994).

¹⁵³ *Authors Guild*, 804 F.3d at 221.

¹⁵⁴ *Harper & Row Publishers, Inc. v. Nation Enters.*, 471 U.S. 539, 565 (1985) (“[A] taking may not be excused merely because it is insubstantial with respect to the *infringing* work.”).

¹⁵⁵ *Id.* at 543.

¹⁵⁶ *Id.*

canceled its piece.¹⁵⁷ Despite the brevity of the excerpt, the Court found that *The Nation's* use weighed against fair use because the excerpt included key revelations about Ford's decision to pardon Nixon that got to the heart of the author's personal reflections.¹⁵⁸

Where the defendant's copying of the original work is necessary to achieve its transformative purpose, the amount and substantiality factor weighs in favor of a finding of fair use.¹⁵⁹ In *Authors Guild*, because Google's use of the original works was to create a search function, it necessarily needed to copy the entire original work, or else its search function would not be able to perform its intended, and transformative, purpose: aid online searchers with determining whether their searched term appears in a book.¹⁶⁰ As a result, even though Google used the entire original work, and thus its use was substantial in amount, the court held that in terms of the search function, the amount and substantiality factor weighed in favor of fair use.¹⁶¹ Moreover, Google's second challenged use, its snippet feature, was not found to be so substantial in amount as to weigh against fair use.¹⁶² Although the plaintiff was able to provide evidence that some searches revealed nearly sixteen percent of their book's text, Google's snippet use is not substantial because the user saw only scattered and fragmented sections of the book's text.¹⁶³ In this way, there was not a significant risk that Google's use would act as a substitute for the original books because it did not coherently communicate the content of the original book.¹⁶⁴

As shown in Exhibit J,¹⁶⁵ the fact that a substantial portion of the infringing work is copied nearly verbatim is powerful evidence of the qualitative importance of the material both to the original creator and to the infringer.¹⁶⁶ Again, GPT-4's dataset is relatively unknown;¹⁶⁷ however, the *Times* alleges that its webpage domain, NYTimes.com, is one of the "'top 15 domains by volume' in the WebText dataset, and is listed as the 5th 'top domain' in the WebText dataset with 333,160

¹⁵⁷ *Id.*

¹⁵⁸ *Id.*

¹⁵⁹ *Authors Guild v. Google, Inc.*, 804 F.3d 202, 221 (2d Cir. 2015).

¹⁶⁰ *Id.*

¹⁶¹ *Id.* at 222–23.

¹⁶² *Id.* at 223.

¹⁶³ *Id.*

¹⁶⁴ *Id.* ("If snippet view could be used to reveal a coherent block amount to 16% of a book, that would raise a very different question beyond the scope of our inquiry.")

¹⁶⁵ See New York Times Complaint, *supra* note 4, at 30–37.

¹⁶⁶ *Harper & Row Publishers, Inc. v. Nation Enters.*, 471 U.S. 539, 565 (1985).

¹⁶⁷ See *supra* text accompanying notes 50–54.

entries.”¹⁶⁸ Even if the use of the *Times*’s articles is unintentional or necessary for training, the model’s verbatim reproduction of paragraphs, OpenAI’s stated aim to use the highest-quality training materials, and the *Times*’s outsized presence in WebText strongly support the *Times*’s claims. After all, if only small fragments or minimal samples had been included, such detailed recitations would be far less likely.¹⁶⁹

4. *The Use’s Effect on the Market Value of the Copyrighted Work*

The last factor is how the defendant’s use of the original work affects “the potential market for or value of the copyrighted work.”¹⁷⁰ This factor examines whether the defendant’s use functions as a market substitute for the original work, thereby causing the copyright owner economic harm by significantly reducing potential revenues due to consumer preference for the defendant’s copy over the original.¹⁷¹ This factor is essential in the fair use analysis because copyright is ultimately a “commercial doctrine” that aims to encourage creativity by ensuring that authors profit from their creations.¹⁷² Moreover, this factor is closely linked with the first factor in the fair use defense analysis (the purpose and character of the use) in that “the more the copying is done to achieve a purpose that differs from the purpose of the original, the less likely it is that the copy will serve as a satisfactory substitute for the original.”¹⁷³

Where a defendant’s use can possibly, probably, or even certainly cause some loss of sales, the loss of sales alone is insufficient to make the defendant’s copy an effective substitute that would weigh against a finding of fair use.¹⁷⁴ In *Authors Guild*, the court found the fourth factor did not weigh against a finding of fair use because Google’s

¹⁶⁸ New York Times Complaint, *supra* note 4, at 25.

¹⁶⁹ *See id.* at 30–37.

¹⁷⁰ *Campbell v. Acuff-Rose Music, Inc.*, 510 U.S. 569, 570 (1994).

¹⁷¹ *Authors Guild v. Google, Inc.*, 804 F.3d 202, 223 (2d Cir. 2015).

¹⁷² *Id.*; *see also* *Harper & Row Publishers, Inc. v. Nation Enters.*, 471 U.S. 539, 566 (1985) (describing the fourth factor as “undoubtedly the single most important element of fair use”).

¹⁷³ *Authors Guild*, 804 F.3d at 223 (citing *Campbell*, 510 U.S. at 591).

¹⁷⁴ *Id.* at 224 (holding that the fourth factor weighs in favor of fair use despite the recognition that the “snippet function can cause *some* loss of sales . . . [Because] the possibility, or even the probability or certainty, of some loss of sales does not suffice to make the copy an effectively competing substitute that would tilt the weighty fourth factor in favor of the rights holder in the original”).

snippet function was so brief, “cumbersome, disjointed, and incomplete” it could hardly be said that the snippet view could act as a significant substitute for the author’s entire book.¹⁷⁵

This factor is likely to weigh against the application of the fair use defense. Here, OpenAI’s use undermines the *Times*’s ability to license content and threatens existing and potential revenue streams, including through content syndication and licensing arrangements. Therefore, unlike the fragmented search facilitation at issue in *Authors Guild*, where the output consisted of isolated, non-substitutive snippets, OpenAI’s ingestion and deployment of the *Times*’s works risks enabling outputs that function as close substitutes for the original articles. ChatGPT’s ability to generate outputs resembling news reporting—sometimes in substantial portions of the original text—means that users could turn to OpenAI’s platform for information rather than consulting original journalism. This dynamic directly undermines both the market for the *Times*’s work and the exclusivity that copyright law is designed to protect. In this respect, OpenAI’s use arguably crosses a line: Rather than serving a distinct transformative function, it imperils the very market in which the *Times* operates, weighing heavily against a finding of fair use under the fourth factor.

V

BROADER LEGAL AND SOCIAL CONSEQUENCES

This Part evaluates the policy implications of a verdict for either OpenAI or *The New York Times*, ultimately concluding that the latter better serves the constitutional purpose of copyright. It first examines the broader policy consequences of a verdict favoring the *Times*, including concerns about algorithmic bias and the limitation of training data to public domain or licensed sources. While these concerns raise important questions about representation and technological development, they do not justify weakening the core protections of copyright. The Part then turns to the implications of a ruling in favor of the *Times*, highlighting how enforcing copyright in the AI context promotes creative sustainability and protects the livelihoods of authors. Although both sides raise compelling concerns, this Part argues that the policy arguments ultimately favor the *Times*, whose position better aligns with the Copyright Clause’s foundational goal: incentivizing

¹⁷⁵ *Id.* at 224–25.

human creativity by granting authors meaningful control over the use of their work.¹⁷⁶

A. The Possibility of Algorithmic Gatekeeping

If the *Times* prevails in its lawsuit and the court mandates the destruction of GPT's training datasets, the implications will be far-reaching. If existing datasets were disallowed, thus forcing AI developers to rebuild, developers would face a critical decision: either rely exclusively on public domain works or selectively license copyrighted materials. Existing datasets are curated through web crawlers, whose objective is to obtain as many materials as possible for the LLM's training, regardless of the source or content.¹⁷⁷ While the training materials are inherently biased by virtue of societal conditions,¹⁷⁸ the risks associated with adding an additional layer of discretionary bias would only exacerbate and amplify that problem.¹⁷⁹ This additional level of discretionary bias is particularly harmful given that "AI systems enjoy an aura of objectivity and accuracy."¹⁸⁰ This appearance of neutrality may result in a reinforcement and normalization of existing prejudices through supposedly factual outputs.

Ultimately, if AI developers relied solely on public domain works, an entirely new set of problems related to bias may be introduced. Copyrighted works from 1929 entered the U.S. public domain on

¹⁷⁶ U.S. CONST. art. I, § 8, cl. 8.

¹⁷⁷ See *Overview of OpenAI Crawlers*, OPENAI PLATFORM, <https://platform.openai.com/docs/bots/> (on file with the Oregon Law Review) (last visited Apr. 14, 2025); Duncan Anderson, *Beyond Data Hoovering: The Nuanced Reality of Training Large Language Models (LLMs)*, MEDIUM (July 19, 2023), <https://medium.com/barnacle-labs/beyond-data-hoovering-the-nuanced-reality-of-training-large-language-models-llms-79aa107c17db> [<https://perma.cc/9A6U-L9Z7>].

¹⁷⁸ Thomas Dethmann & Jannis Spiekermann, *Ethical Use of Training Data: Ensuring Fairness and Data Protection in AI*, LAMARR INST. FOR MACH. LEARNING & A.I. (July 3, 2024), <https://lamarr-institute.org/blog/ai-training-data-bias/> [<https://perma.cc/9UBE-29H8>].

¹⁷⁹ See Uwe Peters, *Algorithmic Political Bias in Artificial Intelligence Systems*, 35 PHIL. & TECH., Mar 2022, at 25, 2 (2022), <https://pmc.ncbi.nlm.nih.gov/articles/PMC8967082/> [<https://perma.cc/K6VF-VCGH>] (Algorithmic bias is defined as "a tendency to not merely neutrally transform or extract information from data but to operate on it in ways that deviate from a normative (moral, statistical, social, etc.) standard such that one kind of individual or group is unfairly privileged over another based on aspects of their social identity.").

¹⁸⁰ *Id.*

January 1, 2025.¹⁸¹ Hence, because the materials in the public domain are from 1929 and before, these materials predominantly skew toward older, Western-centric, and historically privileged perspectives.¹⁸² The historic exclusion of marginalized voices in publishing contribute to the disproportionate representation of Western-centric perspectives.¹⁸³ Undoubtedly, while the publication of literary works in the modern age still portrays certain class-based biases, those biases were inflated in the early twentieth century when “publication was the near-exclusive domain of white males born to a privileged class” and although “exceptions certainly exist[ed], they were relatively rare and special cases.”¹⁸⁴ As a result, AI models trained solely on public domain works risk perpetuating outdated viewpoints while also failing to capture contemporary and diverse discourses. These biases, in turn, affect the inclusivity and objectivity of AI-generated knowledge.

The process of selecting and licensing training data is neither neutral nor purely technical—it is an act of curation that reflects the individual priorities, limitations, and biases of AI developers. The economic considerations underlying AI training data licensing can significantly affect the diversity of the content that is incorporated into models. To be able to use works outside the public domain, AI developers will need to negotiate a license with the relevant copyright owners.¹⁸⁵ There are license fees and transaction costs associated with the negotiation of a license.¹⁸⁶ The license fee is the actual price a company pays to obtain the rights to use the licensed material whereas the transaction costs are the expenses incurred to make the transaction

¹⁸¹ Jennifer Jenkins & James Boyle, *January 1, 2025 Is Public Domain Day: Works from 1929 Are Open to All, as Are Sound Recordings from 1924!*, DUKE L.: CTR. & PROGRAM: CTR. FOR THE STUDY OF THE PUB. DOMAIN, <https://web.law.duke.edu/cspd/publicdomainday/2025/#:~:text=On%20January%201%2C%202025%2C%20thousands,will%20now%20be%20public%20domain> [https://perma.cc/PWL9-HHJJ].

¹⁸² See Cathay Y. N. Smith, *Editing Classic Books: A Threat to the Public Domain?*, 110 VA. L. REV. 1, 5–7 (2024).

¹⁸³ Catherine Tucker, *Potential Socioeconomic Biases of AI Policy*, in IDENTIFYING THE ECONOMIC IMPLICATIONS OF ARTIFICIAL INTELLIGENCE FOR COPYRIGHT POLICY: CONTEXT AND DIRECTION FOR ECONOMIC RESEARCH, U.S. COPYRIGHT OFFICE, 54, 55 (Brent A. Lutes ed., 2025), <https://www.copyright.gov/economic-research/economic-implications-of-ai/Identifying-the-Economic-Implications-of-Artificial-Intelligence-for-Copyright-Policy-FINAL.pdf> [https://perma.cc/9WUW-22AP].

¹⁸⁴ *Id.*

¹⁸⁵ See *Public Domain and Creative Commons: A Guide to Works You Can Use Freely*, UNIV. OF MONT. MAUREEN & MIKE MANSFIELD LIBR. (Jan. 8, 2025, at 16:31 PT), <https://libguides.lib.umt.edu/PublicDomainCC> [https://perma.cc/6Q3K-XN6X].

¹⁸⁶ Richard A. Posner, *Transaction Costs and Antitrust Concerns in the Licensing of Intellectual Property*, 4 J. MARSHALL REV. INTELL. PROP. L. 325, 325 (2005).

happen.¹⁸⁷ Therefore, as more individual owners become involved in the licensing negotiations, transaction costs increase because more labor and money is needed to obtain all the necessary permissions.¹⁸⁸ This economic reality incentivizes developers to negotiate with large, consolidated content providers rather than engaging with a fragmented array of rights holders.¹⁸⁹ This preference for efficiency reinforces the overrepresentation of dominant, corporate-controlled content while sidelining alternative perspectives. For example, corporate news outlets often set the agenda for public discourse,¹⁹⁰ meaning that AI models trained primarily on their content may be more likely to reproduce viewpoints that align with mainstream, profit-driven media rather than those emerging from independent investigative journalism, activist publications, or community-driven storytelling. As a result, AI-generated knowledge may lack the diversity of thought necessary to present a well-rounded, multifaceted perspective on complex issues.

Moreover, given that AI developers are more likely to license works from major corporations, and because those news corporations generally engage in agenda setting,¹⁹¹ there is a potential for political bias. This bias may emerge in two ways: first, through the exclusion of radical or nonmainstream sources from training data; and second, through the prioritization of political works that economically or politically benefit the AI company. Existing algorithms used by website operators and social media platforms, for instance, are often designed to proactively block or remove content deemed politically sensitive, dangerous, or untrustworthy.¹⁹² While filtering certain types of content—such as explicit hate speech—may be widely accepted as necessary, the implementation of policy-directed filtering

¹⁸⁷ *Id.* (e.g., the costs of finding the owner of the material, drafting contracts, paying legal counsel, etc.).

¹⁸⁸ Tucker, *supra* note 183, at 57.

¹⁸⁹ *Id.* at 56–57 (“Reddit . . . has announced a deal to license its users’ postings to developers. However, such forums are poor representations of society at large Reddit is only used by 11 percent of U.S. adults. Of those, around two-thirds are men.” (internal citation omitted)).

¹⁹⁰ See Gary King et al., *How The News Media Activate Public Expression and Influence National Agendas*, 358 *SCIENCE* 776 (2017); Maxwell E. McCombs & Donald L. Shaw, *The Agenda-Setting Function of Mass Media*, 36 *THE PUB. OPINION Q.* 176, 177 (1972), <https://www.jstor.org/stable/2747787> (“[T]he press ‘may not be successful much of the time in telling people what to think, but it is stunningly successful in telling its readers what to think about.’”).

¹⁹¹ King, *supra* note 190.

¹⁹² Peters, *supra* note 179, at 5.

algorithms may “violate legitimate user expectations of policy-neutrality” and contribute to ideological discrimination by disproportionately suppressing certain perspectives.¹⁹³ Taken to the extreme, such practices may hinder the free marketplace of ideas, thereby restricting public debate, limiting open dialogue, and enabling de facto censorship.¹⁹⁴ Additionally, while it is implausible—if not impossible—to capture every conceivable political work from across the ideological spectrum, the discretion AI developers exercise in determining which political materials are *relevant* enough to warrant the financial and logistical effort of licensing is problematic. This selection process effectively functions as a form of self-censorship, wherein AI-generated content disproportionately aligns with dominant ideological frameworks rather than presenting a truly comprehensive understanding of political issues.¹⁹⁵ The exclusion of diverse political viewpoints from AI training datasets ultimately risks reinforcing existing power structures and narrowing the range of perspectives available in public discourse.

Similarly, AI developers may prioritize obtaining licenses from established publishing institutions over less conventional sources. This preference is, in part, justified by the perception that such institutions transmit objective, factual information curated by credentialed professionals.¹⁹⁶ However, this assumption overlooks the historical association of mainstream media with political and economic interests.¹⁹⁷ Since at least the nineteenth century, newspapers and other major media outlets have frequently aligned themselves with specific political parties and economic power structures, shaping public discourse in ways that reflect their affiliations rather than practicing purely neutral reporting.¹⁹⁸ The historical bias in media is further compounded by the fact that dominant publishing institutions have traditionally been owned and operated by elite groups with vested

¹⁹³ *Id.*

¹⁹⁴ G. Michael Parsons, *Fighting for Attention: Democracy, Free Speech, and the Marketplace of Ideas*, 104 MINN. L. REV. 2157, 2158–59 (2020).

¹⁹⁵ Anderson, *supra* note 177 (“We want our model to understand lots of forms and styles of language, so it’s important that we include a very wide variety of language examples in our training data. If a model has never seen two humans arguing and insulting each other, it cannot understand such. Similarly, a model that’s not seen street language and swearing will struggle to understand it. That’s why the training data for an LLM needs to include data that might raise an eyebrow or two. If it doesn’t, our model will be naive.”).

¹⁹⁶ Bruce Thornton, *A Brief History of Media Bias*, HOOVER INST. (June 12, 2023), <https://www.hoover.org/research/brief-history-media-bias> [<https://perma.cc/WB6S-9XGE>].

¹⁹⁷ *Id.*

¹⁹⁸ *Id.*

interests in maintaining prevailing power structures.¹⁹⁹ Meanwhile, alternative and grassroots publications—often serving as critical counterpoints to mainstream narratives—have historically struggled to achieve the same levels of influence and legitimacy.²⁰⁰ Therefore, if AI models disproportionately train on sources that have historically filtered public knowledge through elite-controlled lenses, they may inadvertently reinforce systemic exclusions and fail to provide a truly comprehensive or equitable representation of knowledge.

Ultimately, the biases embedded in licensing decisions have significant consequences for the reliability, inclusivity, and fairness of AI-generated knowledge. When certain voices and perspectives are excluded from training datasets, AI models generate content that reflects and reinforces these exclusions. This can result in skewed, incomplete, and misleading AI-generated information; marginalization of underrepresented communities; reinforcement of existing media and academic biases; and overarching fairness and accountability concerns. Hence, bias in AI training data is not limited to the content that is excluded; it also exists in the choices made during the licensing process. When AI developers prioritize certain sources over others based on economic, political, or systemic biases, they shape the perspectives that AI models ultimately produce. Addressing these biases requires a commitment to transparency, inclusivity, and ethical responsibility in AI development.

B. The Antithesis of the Copyright Clause

If OpenAI successfully establishes a fair use defense for its GPT training, the resulting precedent could undermine the core purposes of copyright law—encouraging creativity and ensuring economic incentives for authors—by allowing AI companies to exploit copyrighted works without paying compensation.²⁰¹ By granting

¹⁹⁹ Kate Vinton, *These 15 Billionaires Own America's News Media Companies*, FORBES (Jun. 1, 2016, at 14:26 ET), <https://www.forbes.com/sites/katevinton/2016/06/01/these-15-billionaires-own-americas-news-media-companies/> [<https://perma.cc/2VHJ-RB8D>] (The Murdoch family “controls 120 newspapers across five countries” and the Cox family owns “more than a dozen non-daily publications, 14 broadcast television stations, one local cable channel and 59 radio stations.”).

²⁰⁰ Karoline Andrea Ihlebæk et al., *Understanding Alternative News Media and Its Contribution to Diversity*, 10 DIGIT. JOURNALISM, 1267, 1268–69 (2022), <https://www.tandfonline.com/doi/full/10.1080/21670811.2022.2134165?scroll=top&needAccess=true#abstract> [<https://perma.cc/8HLJ-DGFR>].

²⁰¹ U.S. CONST. art. I, § 8, cl. 8.

creators exclusive rights over their works, copyright law fosters investment in intellectual labor, enabling authors, journalists, and publishers to sustain their livelihoods.²⁰² However, if courts accept a broad fair use defense for AI training, this foundational purpose could be undermined. AI companies, such as OpenAI, profit from models trained on copyrighted materials without permission or compensation.²⁰³ This practice threatens the economic incentives that sustain creative industries and commercially exploits the labor of authors. If AI-generated content can directly compete with human-created works, demand for original literature, journalism, and other forms of creative expression will decline, leading to fewer new works and a diminished cultural landscape.²⁰⁴

By providing temporary monopolies over creative works, copyright law ensures that authors are rewarded for their contributions, thereby encouraging the continued production of literature, journalism, music, and other forms of expression.²⁰⁵ Without strong copyright protections, authors may lose the motivation to create new works, knowing that their content could be freely exploited without compensation.²⁰⁶ The exclusivity granted by copyright law fosters a balanced and sustainable knowledge environment—one in which creators can earn a livelihood while contributing to cultural and intellectual development.²⁰⁷ A broad fair use defense for AI training undermines this principle by enabling AI firms to appropriate copyrighted materials without permission, stripping authors of their economic rights, and weakening

²⁰² Mei-Ian Stark, *5 Ways Copyright Laws Encourage Personal Expression and Creativity*, U.S. CHAMBER OF COM. (Apr. 25, 2022), <https://www.uschamber.com/intellectual-property/five-ways-copyright-laws-encourage-personal-expression-and-creativity> [<https://perma.cc/DB98-ENP8>].

²⁰³ Milmo, *supra* note 14.

²⁰⁴ See New York Times Complaint, *supra* note 4, at 14.

²⁰⁵ *Copyright Basics*, U.S. PAT. & TRADEMARK OFF., <https://www.uspto.gov/ip-policy/copyright-policy/copyright-basics#:~:text=The%20primary%20purpose%20behind%20copyright,them%20available%20in%20the%20marketplace> [<https://perma.cc/ZV9U-L5TD>] (“By granting authors the exclusive right to authorize certain uses of their works, copyright provides economic incentives to create new works and to make them available in the marketplace”).

²⁰⁶ See Celeste Shen, *Fair Use, Licensing, and Authors’ Rights in the Age of Generative AI*, 22 NW. J. TECH. & INTELL. PROP. 157, 177 (2024); see also Harry Jiang et al., *AI Art and Its Impact on Artists*, AIES 363, 368 (2023), <https://doi.org/10.1145/3600211.3604681> (on file with the Oregon Law Review) (In the context of visual arts, “students who foresee image generators replacing artists have become demoralized and dissuaded from honing their craft and developing their style.”).

²⁰⁷ See *Andy Warhol Found. for the Visual Arts, Inc. v. Goldsmith*, 598 U.S. 508, 526 (2023).

the incentives necessary for creative production. Additionally, large technology firms, which already control vast amounts of digital content,²⁰⁸ stand to further consolidate their power by freely using copyrighted works without compensating creators. This creates an imbalance, in which AI companies profit from content they did not create, while smaller publishers and independent authors struggle to compete.

Moreover, copyright protections ensure that authors and publishers are compensated through book sales, licensing fees, and royalties.²⁰⁹ However, if AI models are allowed to freely train on copyrighted materials, they can generate content that competes directly with human-created works.²¹⁰ This could significantly reduce demand for original literature and journalism.²¹¹ Already, some authors and media outlets have reported declining revenues as AI-generated summaries and articles divert traffic from their platforms.²¹² Without financial incentives, many writers may find it unsustainable to continue producing new works, leading to a decline in creative output.

Ultimately, when authors see their works freely exploited by AI companies without consent, they may lose the motivation to create.²¹³ Fear that original works will be absorbed into AI training datasets without compensation or recognition discourages creative risk-taking and innovation.²¹⁴ AI-generated works are inherently derivative, as they rely on existing copyrighted materials to produce new content.²¹⁵ Consequently, if AI-generated literature and journalism become the widespread dominant form of media, the

²⁰⁸ Shivam Rapjot, *How Does Big Tech Dominate the Digital World?*, GEOSTRATA (Sept. 30, 2024), <https://www.thegeostrata.com/post/big-tech-tightens-its-grip-how-companies-like-apple-and-google-dominate-the-digital-world> [https://perma.cc/XK39-3MJ3] (For example, “Meta Inc. is a leader in social media and communication . . . with over 3 billion active users.”).

²⁰⁹ 17 U.S.C. § 106.

²¹⁰ See New York Times Complaint, *supra* note 4, at 14.

²¹¹ See *id.* at 14–15.

²¹² Wendy Lee, *Get Paid or Sue? How the News Business Is Combating the Threat of AI*, L.A. TIMES (July 24, 2024, at 03:00 PT), <https://www.latimes.com/entertainment-arts/business/story/2024-07-24/how-ai-is-disrupting-the-journalism-industry> [https://perma.cc/72XK-MAZ6].

²¹³ *Copyright Basics*, *supra* note 205.

²¹⁴ See Jiang et al., *supra* note 206, at 368 (“[B]oth new and current artists are becoming increasingly reluctant to share their works and perspectives, in an attempt to protect themselves from the mass scraping and training of their life’s works.”); see also *Copyright Basics*, *supra* note 205.

²¹⁵ See *supra* text accompanying notes 31–33.

diversity of artistic expression may suffer. Instead of fostering new ideas and voices, AI-generated content could reinforce dominant narratives, leading to a homogenization of creative output.²¹⁶

Further, artists may hesitate to share their work and mentor aspiring creators, which not only stifles individual artistic growth but also weakens the collective creativity of humanity.²¹⁷ Artistic innovation has long relied on a rich tradition of learning, adaptation, and reinterpretation, where new generations of creators draw inspiration from those who came before them, building upon past techniques to develop novel expressions.²¹⁸ When this cycle is disrupted by the fear of exploitation, fewer artists may engage in this exchange of knowledge, resulting in a narrower landscape of creativity. This phenomenon mirrors a troubling feedback loop emerging in AI development, particularly in the training of LLMs.²¹⁹ When AI models are trained on the outputs of previous models rather than on a diverse corpus of human-created content, they risk generating increasingly homogenized and inaccurate results.²²⁰ Each successive generation becomes more detached from the originality, complexity, and novelty of human expression, and instead amplifies the patterns, biases, and inherent limitations from earlier iterations of AI models.²²¹ This self-referential process resembles an ouroboros—a serpent consuming its own tail—where AI-generated content continuously feeds into new models, creating a cycle in which nothing truly original is produced.²²²

²¹⁶ See Anil R. Doshi & Oliver P. Hauser, *Generative AI Enhances Individual Creativity but Reduces the Collective Diversity of Novel Content*, NAT'L LIBR. MED. (July 12, 2024), <https://pmc.ncbi.nlm.nih.gov/articles/PMC11244532/> [<https://perma.cc/QXD5-XUEQ>].

²¹⁷ Jiang et al., *supra* note 206, at 368.

²¹⁸ See, e.g., Daniela Lehner, *The Heroine's/Hero's Journey—A Call for Transformation? Transformative Learning, Archetypal Patterns, and Embodied Knowing/Learning*, 20 J. TRANSFORMATIVE EDUC. 88, 89 (2022), <https://journals.sagepub.com/doi/10.1177/15413446211007089?icid=int.sj-abstract.citing-articles.21> (on file with the Oregon Law Review) (explaining the “Heroine/Hero’s Journey” is a common storytelling trope that can be traced through “myths and stories from different times and places”).

²¹⁹ Iliia Shumailov et al., *AI Models Collapse When Trained on Recursively Generated Data*, 631 NATURE, 755, 755 (2024), <https://www.nature.com/articles/s41586-024-07566-y> [<https://perma.cc/V3PF-DR8F>] (“Model collapse is a degenerative process affecting generations of learned generative models, in which the data they generate end up polluting the training set of the next generation. Being trained on polluted data, they then mis-perceive reality.”).

²²⁰ Keyur Ramoliya, *Model Collapse in AI*, MEDIUM (Aug. 2, 2024), <https://medium.com/the-deephub/model-collapse-in-ai-813418fd8516> [<https://perma.cc/7VDU-59SH>].

²²¹ *Id.* (As the feedback loop continues, “each new AI model learns more from AI-created content and less from human-created content . . .” which “can lead to a gradual degradation in the quality and diversity of AI-generated content.”).

²²² Jiang et al., *supra* note 206, at 368.

If humanity comes to rely exclusively on AI-generated works for “the words we read, the art we see,” and the music we hear, our cultural output may become an echo chamber of past creations rather than a space for genuine innovation.²²³ Instead of pushing artistic and intellectual boundaries, we risk a future where creativity is reduced to an automated regurgitation of the *familiar*, rather than nurturing the expansion of dynamic creative expression.

C. More Than Just a © Symbol

A verdict favoring *The New York Times* aligns squarely with the foundational purpose of the Copyright Clause of the U.S. Constitution: “[t]o promote the Progress of Science and useful Arts.”²²⁴ The Copyright Clause, therefore, reflects a deliberate constitutional choice to empower creators with exclusive rights as a means to incentivize the continuous production of original works, thereby fostering a flourishing, innovative society.²²⁵ The quid pro quo nature of copyright law—providing a time-limited monopoly in exchange for public benefit—depends on the enforceability of copyright protections; even, and especially, in the face of transformative technological shifts like generative AI.²²⁶ If AI developers are permitted to harvest and commercially exploit copyrighted journalism without consent, compensation, or credit, the constitutional calculus underpinning copyright unravels.²²⁷ Creators lose the incentive to invest the time, labor, and resources needed for high-quality content, and it is the public who ultimately suffers from a diminished cultural and informational landscape.²²⁸

Although enforcing copyright in the AI training context could contribute to algorithmic bias or selective representation,²²⁹ these concerns cannot justify eroding the fundamental principles of copyright law. The Copyright Clause was never meant to guarantee equal data

²²³ *Id.*

²²⁴ U.S. CONST. art. I, § 8, cl. 8.

²²⁵ *Artl.S8.C8.1 Overview of Congress’s Power Over Intellectual Property*, CONST. ANNOTATED, https://constitution.congress.gov/browse/essay/artl-S8-C8-1/ALDE_00013060/ (on file with the Oregon Law Review) (last visited Oct. 6, 2025).

²²⁶ *Id.* (“Without legal protection, competitors could freely copy such creations, denying the original creators the ability to recoup their investments in time and effort, reducing the incentive to create in the first place.”).

²²⁷ *Id.*

²²⁸ *Id.*

²²⁹ See *supra* text accompanying notes 175–98.

access to private companies building profitable AI systems; rather, it is a guarantee of exclusive rights for authors, inventors, and other creators.²³⁰ Moreover, the risks of bias or content gatekeeping arise not from respecting those rights, but from deeper issues like media monopolization and technology governance.²³¹ Consequently, honoring the *Times*'s exclusive right to its work in the AI context affirms a constitutional commitment to incentivize human creativity and compensate the labor of creatives, irrespective of the potential policy consequences.

CONCLUSION

The *Times*'s lawsuit against OpenAI raises fundamental questions about the intersection of copyright law, fair use, and the ethical responsibilities of AI developers. At the heart of this dispute lies the challenge of balancing AI innovation with the legal and economic protections afforded to authors, journalists, and other creators. Whereas OpenAI asserts that its use of copyrighted materials for training LLMs constitutes fair use,²³² the *Times* contends that such practices undermine the economic incentives that sustain creative industries.²³³

A fair use analysis of OpenAI's practices suggest that although AI training may satisfy certain fair use factors, the commercial nature of OpenAI's operations and the potential market harm to original creators' work weigh against a broad fair use defense.²³⁴ The nature of AI-generated content—derivative yet transformative—complicates fair use jurisprudence, highlighting the need for courts to consider the unique attributes of machine learning models in assessing copyright claims.²³⁵

Beyond the legal arguments, this case carries significant implications for the broader media and technology landscape. If the *Times* prevails and courts mandate the destruction of OpenAI's training datasets, AI developers may be forced to reconstruct their models using either public domain works or selectively licensed materials.²³⁶ This shift risks introducing new biases into AI-generated content, as reliance on public domain works skews toward older, historically privileged

²³⁰ See U.S. CONST. art. I, § 8, cl. 8.

²³¹ See *supra* text accompanying notes 187–98.

²³² Defendant's Memorandum in Support of Motion to Dismiss, *supra* note 9, at 3.

²³³ New York Times Complaint, *supra* note 4, at 14.

²³⁴ See *supra* Section IV.B.

²³⁵ *Id.*

²³⁶ See *supra* text accompanying notes 175–98.

perspectives, while selective licensing favors dominant, corporate-controlled narratives.²³⁷ The economic and political consequences of algorithmic gatekeeping could further entrench existing power structures, limiting the diversity of perspectives represented in AI-generated knowledge.²³⁸

Conversely, if OpenAI's fair use defense is upheld, it could weaken copyright protections that ensure economic incentives for creators. Allowing AI companies to profit from models trained using copyrighted works without compensation or credit may undermine the sustainability of journalism, literature, and other creative industries.²³⁹ Without adequate safeguards, AI-generated content could displace human authors, thereby reducing the financial viability of creative professions and diminishing the diversity of cultural and intellectual discourse more broadly.²⁴⁰

Ultimately, this case underscores the urgent need for legal frameworks that address the evolving relationship between AI and copyright law. This tension between technological advancement and the protection of creators' rights is not a new phenomenon.²⁴¹ Throughout history, groundbreaking innovations have disrupted existing creative industries, promoting legal and ethical debates about ownership, authorship, and fair use. For example, nineteenth-century painters and illustrators resisted the advent of the photograph, as they feared that mechanically produced images would devalue traditional artistic skills.²⁴² Similarly, the rise of the internet and digital media in the late twentieth and early twenty-first centuries led to widespread concerns over copyright infringement, as file-sharing platforms and

²³⁷ *Id.*

²³⁸ *Id.*

²³⁹ See *supra* Section V.B.

²⁴⁰ *Id.*

²⁴¹ Eva Silva, *How Photography Pioneered a New Understanding of Art*, COLLECTOR (June 4, 2022), <https://www.thecollector.com/how-photography-transformed-art/> [<https://perma.cc/CVB5-79NE>]. See generally Jean K. Chalaby, *The Streaming Industry and the Platform Economy: An Analysis*, 46 MEDIA CULTURE & SOC'Y 552, 554, <https://journals.sagepub.com/doi/10.1177/01634437231210439#:~:text=SVoD%20services%20can%20acquire%20the,Turton%20and%20Opie%2C%202019> (on file with the Oregon Law Review); *The Evolution and Impact of Spotify on the Music Industry*, MEDIUM: DIGITAL MIRAI (Jul. 15, 2024), <https://digitalmirai.medium.com/the-evolution-and-impact-of-spotify-on-the-music-industry-49a140cf9735> [<https://perma.cc/3MH2-4WWF>] (“Spotify has democratized music distribution, allowing independent artists to reach global audiences without the need for traditional record deals.”).

²⁴² Silva, *supra* note 241.

streaming services reshaped the ways creative works were distributed and monetized.²⁴³ Most recently, the popular music platform Spotify disrupted the music industry's strict control over music distribution by making music widely accessible through a digital platform.²⁴⁴ In all cases, legal frameworks and industry practices ultimately evolved to balance technological progress with the rights of creators. The debate over AI-generated content is just the latest iteration of this ongoing struggle, raising a fundamental question about how to foster innovation without undermining the very human creativity that fuels it.

²⁴³ See generally Chalaby, *supra* note 241.

²⁴⁴ *The Evolution and Impact of Spotify on the Music Industry*, *supra* note 241.