

Clarkson Law Firm, P.C. | 22525 Pacific Coast Highway, Malibu, CA 90265 | P: (213) 788-4050 | F: (213) 788-4070 | clarksonlawfirm.com

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

**CLARKSON LAW FIRM, P.C.**  
Ryan J. Clarkson (CA SBN 257074)  
*rclarkson@clarksonlawfirm.com*  
Yana Hart (CA SBN 306499)  
*yhart@clarksonlawfirm.com*  
Tiara Avanness (CA SBN 343928)  
*tavaness@clarksonlawfirm.com*  
22525 Pacific Coast Highway  
Malibu, CA 90265  
Tel: (213) 788-4050

**CLARKSON LAW FIRM, P.C.**  
Tracey Cowan (CA SBN 250053)  
*tcowan@clarksonlawfirm.com*  
95 3rd St., 2nd Floor  
San Francisco, CA 94103  
Tel: (213) 788-4050

*Counsel for Plaintiff and the Proposed Class*

**UNITED STATES DISTRICT COURT  
NORTHERN DISTRICT OF CALIFORNIA**

PLAINTIFF JILL LEOVY, individually, and on  
behalf of all others similarly situated,

Plaintiff,

vs.

GOOGLE LLC,

Defendant.

Case No. 3:23-cv-3440-AMO

**SECOND AMENDED CLASS ACTION  
COMPLAINT**

- 1. DIRECT COPYRIGHT  
INFRINGEMENT

**DEMAND FOR JURY TRIAL**

**TABLE OF CONTENTS**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

INTRODUCTION .....1

PARTIES .....4

JURISDICTION AND VENUE.....5

FACTUAL BACKGROUND .....6

I. GOOGLE’S DEVELOPMENT OF GOOGLE GEMINI.....6

    A. Google’s AI Product Development Depends on Stolen Vast Troves of  
    Copyrighted Data, Including Plaintiff’s Book .....10

    B. Google’s Revised Privacy Policy Purports to Give it “Permission” to Take  
    Anything Shared Online to Train and Improve Its AI Products, Including  
    Copyrighted Information ..... 11

    C. Google Uses This Stolen Data to Profit by the Billions .....14

    D. Creators are Outraged by Google’s Theft-Based Training Model .....15

    E. The Public is Outraged by the Lack of Respect for Autonomy in the Copyright  
    Space, and AI Developments Writ Large.....19

    F. Online News and Media Businesses are Taking Action Against Google’s  
    Unauthorized Infringement.....20

II. DEFENDANT’S CONDUCT VIOLATES ESTABLISHED COPYRIGHT LAWS .....21

CLASS ALLEGATIONS .....23

COUNT ONE – DIRECT COPYRIGHT INFRINGEMENT.....27

PRAYER FOR RELIEF .....29

JURY TRIAL DEMANDED.....30

Clarkson Law Firm, P.C. | 22525 Pacific Coast Highway, Malibu, CA 90265 | P: (213) 788-4050 | F: (213) 788-4070 | clarksonlawfirm.com

1 Plaintiff Jill Leovy (“**Plaintiff**”), individually and on behalf of all others similarly situated,  
 2 brings this action against Defendant Google, LLC (“**Defendant**” or “**Google**”). Plaintiff’s  
 3 allegations are based upon personal knowledge as to herself and her own acts, and upon information  
 4 and belief as to all other matters.

### 5 INTRODUCTION

6 1. The Constitution and the Copyright Act recognize the critical importance of  
 7 protecting the authors and creators’ exclusive rights over their works. The legal protection of  
 8 copyrighted materials is intended to nourish and encourage innovation and creativity. As the United  
 9 States Supreme Court has declared: “The immediate effect of our copyright law is to secure a fair  
 10 return for an ‘author’s creative labor. But the ultimate aim is, by this incentive to stimulate artistic  
 11 creativity for the general public good.”<sup>1</sup> Knowing that their works are protected, creators are more  
 12 likely to invest time and effort in creating new works, leading to a richer and more diverse literary  
 13 landscape.

14 2. Google, however, has elected to disregard the Constitution and the Copyright Act, and  
 15 grant itself a license to steal copyright protected works created by millions of Americans to build,  
 16 train, and commercialize its Artificial Intelligence (“AI”) Products like “Gemini” (previously  
 17 known as “Bard”), the chatbot Google released to compete with OpenAI’s “ChatGPT.” For years,  
 18 Google secretly harvested a massive quantity of pirated and copyrighted works, including a trove  
 19 of books, articles, images, photographs, and millions of other protected works. Google used the  
 20 protected works to create its AI Products, and by doing so, increased its market share by billions of  
 21 dollars, unjustly profiting off the infringed intellectual property. Specifically, Google posted that it  
 22 achieved revenue of \$80.5 billion dollars in the first quarter of 2024—which reflects a 15% increase  
 23 from the first quarter of 2023, when it initially released Gemini (originally introduced as Bard).<sup>2</sup>

24  
 25  
 26 

---

<sup>1</sup> *Sony Corp. of Am. v. Universal Studios, Inc.*, 464 U.S. 417, 432, 78 L. Ed. 2d 574, 104 S. Ct. 774 (1984).

27 <sup>2</sup> Nick robins-Early, *Alphabet Hails Once in a Lifetime AI Opportunity as Revenue Rises*, THE  
 28 GUARDIAN (April 25, 2024), <https://www.theguardian.com/technology/2024/apr/25/google-revenue-quarter-one>.

1           3.     Generative AI Products like Gemini are designed to understand and generate human  
2 language. They are characterized by their size and complexity, and are able to write human-like  
3 responses, articles, books, and other works, mimicking the expressive works on which they were  
4 built.

5           4.     Creative and expressive works are critical to the AI training process because this is  
6 how products like Gemini learn to “create” works. This mass theft of personal and copyrighted  
7 information has stunned internet users around the world.

8           5.     The FTC issued a stern warning to the AI industry in May 2023 regarding this sudden  
9 sprint to collect as much training data as they can find: “Machine learning is no excuse to break the  
10 law . . . The data you use to improve your algorithms must be lawfully collected . . . companies  
11 would do well to heed this lesson.”<sup>3</sup>

12          6.     Rather than heed the FTC’s warning and stop its years-long theft of data, Google  
13 elected to copy and download the works of writers and other creators, without compensation, to  
14 build and train its AI Products like Gemini, capable of performing now the same type of work for  
15 which these writers and authors would be paid.

16          7.     Defendant’s billion-dollar AI Products’ success was possible only because it copied,  
17 downloaded, and digested the protected copyrightable works contained in billions of actual texts  
18 across the internet – all without paying a dime to creators.

19          8.     Defendant’s excuse to the massive infringement is that it believes that everything on  
20 the internet is a fair game for Google to take for private gain and commercial use, as it announced  
21 to the public through the updated version of its online privacy policy from July 2023. However,  
22 that’s not the law. The law requires Google to obtain the creators’ consent to reproduce (by copying,  
23 downloading, and storing the works accessed from the online databases or websites) and pay fair  
24 value for such use.

25

---

26           <sup>3</sup> *Statement of Commissioner Alvaro M. Bedoya Joined by Chair Lina M. Khan and Commissioner*  
27 *Rebecca Kelly Slaughter*, FEDERAL TRADE COMMISSION (MAY 31, 2023), available at  
28 <https://www.ftc.gov/legal-library/browse/cases-proceedings/public-statements/statement-commissioner-alvaro-m-bedoya-joined-chair-lina-m-khan-commissioner-rebecca-kelly-slaughter-0>.

1           9. Authors and other creators are outraged to learn that Google has been using their  
2 works to train its AI Products. This reaction is understandable, given Google’s history of  
3 encroaching on intellectual and other data property rights. Its claims of ownership over all internet  
4 content, including copyrighted works, strikes many as bold and insidious—because it is.

5           10. Google has since invited the world to engage in “dialogue” about what data collection  
6 and protection efforts should look like in the new era of AI, while continuing to steal and infringe  
7 the works of authors to improve and expand its AI Products. That invited a backlash of its own,  
8 naturally, as a classic case of too little too late. One commentator aptly translated Google’s  
9 “invitation” into the truth: “Now that we’ve already trained our LLMs on all your proprietary and  
10 copyrighted content, we will finally start thinking about giving you a way to opt out of any of your  
11 future content for being used to make us rich.”<sup>4</sup>

12           11. Google had options other than to steal copyrighted information. Google could have  
13 paid to license the copyrighted material it stole, so as not to infringe on the owner’s exclusive rights.  
14 The legal acquisition of such material depends on consent and consideration.

15           12. There are companies that specialize in curating and selling datasets for AI training  
16 purposes that contain information obtained with the *express consent* of the content creators or  
17 subjects of the personal or copyrighted information. Using these datasets might be more expensive  
18 than stealing, but using this data has one critical advantage: it is legal. Against this backdrop,  
19 Google’s decision to instead take copyrighted material without notice, consent, or fair compensation  
20 not only violates the individual rights of millions, but also gives Google an unfair advantage over  
21 smaller competitors who lawfully license copyrighted material.

22           13. As part of its theft, Google illegally accessed restricted, subscription-based websites  
23 to take the content of millions without permission and infringed on millions of materials explicitly  
24 protected by copyright, including previously stolen property from websites known for pirated  
25 collections of books and other creative works. Without this mass theft of copyrighted information  
26 belonging to real people, many of Google’s AI products, including Gemini, would not exist.

---

27           <sup>4</sup> Matt G. Southern, *Google Calls For Public Discussion On AI Use Of Web Content*, SEARCH  
28 ENGINE JOURNAL (July 7, 2023), <https://www.searchenginejournal.com/google-calls-for-public-discussion-on-ai-use-of-web-content/491053/> (quoting Barry Adams, Twitter/X, since deleted).

1 Defendant, on a mass scale, unlawfully reproduced valuable intellectual property of creators without  
2 paying a penny for the use of these creative works to build its AI Products. Defendant must be  
3 enjoined from these ongoing violations of copyright laws. It must also be ordered to either delete  
4 the data already obtained illegally or obtain licenses to use the copyrighted material it has stolen.

## 5 PARTIES

### 6 Plaintiff Jill Leovy (“Plaintiff Leovy”)

7 14. Plaintiff Leovy is a New York Times best-selling author and investigative journalist  
8 residing in the State of Texas.

9 15. Plaintiff Leovy owns the registered copyright in this book, which includes customary  
10 copyright-management information including the name of the author and the year of publication  
11 (2015). The registered copyright owned by Plaintiff Leovy is included as **Exhibit A**. Defendant  
12 misappropriated Plaintiff Leovy’s award-winning non-fiction book called *Ghettoside: A True Story*  
13 *of Murder in America*, by illegally copying the book in full without her knowledge or consent to  
14 train “Gemini” and Google’s other AI Products.

15 16. The copyrighted work that Defendant misappropriated and otherwise infringed  
16 reflects over a decade of Plaintiff Leovy’s investigative journalism and work, including novel  
17 insights on a topic few have researched and written about in as much detail. As a result of  
18 Defendant’s large-scale theft of copyrighted materials, Plaintiff Leovy was never paid a penny for  
19 Google’s unauthorized reproduction of the book, on which it developed its multi-billion-dollar  
20 product.

21 17. Defendant’s infringement thus deprives creators like Leovy of deserved royalties, and  
22 undermines their financial stability and incentive to create new works. Absent the relief sought in  
23 this Action, Plaintiff Leovy and millions of authors and creators like her presently have no ability  
24 to demand Google to stop copying/reproducing their works, and/or provide fair compensation for  
25 the use of their works.

### 26 Defendant

27 18. **Defendant Google LLC** is headquartered in Mountain View, California. It was  
28 founded in 1998 as an American search engine company. Google’s search business now amounts

1 to \$149 billion, with over 85 percent market share in the global desktop search engine market  
 2 worldwide. In 2015, as part of its corporate restructuring, Google LLC became a subsidiary of its  
 3 newly formed parent company, Alphabet, Inc. Google LLC is currently one of the world’s largest  
 4 for-profit tech companies, specializing in internet related services and products with a special  
 5 emphasis on “web-based search and display advertising tools, search engine, cloud computing,  
 6 software, and hardware.”<sup>5</sup>

7       19. **Agents and Co-Conspirators.** Defendant’s unlawful acts were authorized, ordered,  
 8 and performed by Defendant’s respective officers, agents, employees, representatives, while  
 9 actively engaged in the management, direction, and control of Defendant’s businesses and affairs.  
 10 Defendant’s agents operated under explicit and apparent authority of its principals. Each Defendant,  
 11 and its subsidiaries, affiliates, and agents operated as a single unified entity.

#### 12 JURISDICTION AND VENUE

13       20. This Court has subject matter jurisdiction under 28 U.S.C. § 1331 because this case  
 14 arises under the Copyright Act, 17 U.S.C. § 501.

15       21. Pursuant to 28 U.S.C. § 1391, this Court is the proper venue for this action because a  
 16 substantial part of the events, omissions, and acts giving rise to the claims herein occurred in this  
 17 District: Defendant Google LLC is headquartered in this District, Defendant gains significant  
 18 revenue and profits from doing business in this District, Class Members affected by Defendant’s  
 19 copyright infringement reside in this District, and Defendant employs numerous people in this  
 20 District—a number of whom work specifically on making decisions regarding the handling and use  
 21 of copyrighted materials. Defendant has transacted business, maintained substantial contacts, and/or  
 22 committed overt acts in furtherance of the illegal scheme and conspiracy throughout the United  
 23 States, including in this District. Defendant’s conduct had the intended and foreseeable effect of  
 24 causing injury to persons residing in, located in, or doing business throughout the United States,  
 25 including in this District.

26       22. The Court has general personal jurisdiction over Defendant, because Defendant is

27 \_\_\_\_\_  
 28 <sup>5</sup> *Google LLC*, BLOOMBERG,  
<https://www.bloomberg.com/profile/company/8888000D:US#xj4y7vzkg> (last visited June 27,  
 2024).



1 headquartered in California and makes decisions concerning the Product(s) and the use of use of  
2 copyrighted materials from California. Defendant also advertises and solicits business in California.

### 3 **FACTUAL BACKGROUND**

#### 4 **I. GOOGLE’S DEVELOPMENT OF GOOGLE GEMINI.**

5 23. Defendant’s AI Product, Google Gemini, operates as an advanced language model,  
6 capable swaths of information in response to users’ questions and prompts.<sup>6</sup> Its user interface is  
7 presented as “a dialogue box where users type in their queries.”<sup>7</sup> When users enter their questions  
8 in to Gemini’s dialogue box, Gemini uses the treasure trove of information it has been trained on  
9 and information available on the internet to provide users with responses—often admittedly  
10 plagiarizing the writing of others.<sup>8</sup> Gemini is able to respond to users not only with text-based  
11 answers, but also via image-based answers, adding another function to its capabilities.<sup>9</sup>

12 24. Gemini was initially built on the LaMDA LLM.<sup>10</sup> Google has since transitioned  
13 Gemini to PaLM 2,<sup>11</sup> a LLM trained on 3.6 trillion tokens (strings of words), more powerful than  
14 any existing model.<sup>12</sup> Due to its vast training data, Gemini can also replicate and mimic the human  
15 writing – reproducing the works of the artists, authors, and creators on whose content it was trained  
16 on. The end result is that Gemini is not only built on the works of millions of creators and authors,  
17

18 <sup>6</sup> Andy Patrizio, *Google Bard*, TECHTARGET,  
19 <https://www.techtaraget.com/searchenterpriseai/definition/Google-Bard> (last visited June 27,  
20 2024).

21 <sup>7</sup> Ben Wodecki, *Google Unveils Bard: Its Version of ChatGPT*, AI BUS. (Feb. 7, 2023),  
22 <https://aibusiness.com/google/google-unveils-bard-its-version-of-chatgpt>.

23 <sup>8</sup> See Avram Piltch, *Google Bard Plagiarized Our Article, Then Apologized When Caught*, TOM’S  
24 HARDWARE (March 23, 2023), [https://www.tomshardware.com/news/google-bard-plagiarizing-  
25 article](https://www.tomshardware.com/news/google-bard-plagiarizing-article). The author of this article questioned Google Gemini, at the time Google Bard, about  
26 computer processors, and Gemini provided an answer that was word for word taken form a Tom’s  
27 Hardware article. *Id.* Then, when asked if Gemini had just plagiarized the article, Gemini  
28 responded, “yes what I did was a form of plagiarism.” *Id.*

<sup>9</sup> Sabrina Ortiz, *What is Google Bard? Here’s Everything You Need to Know*, ZDNET (June 1,  
2023), <https://www.zdnet.com/article/what-is-google-bard-heres-everything-you-need-to-know/>.

<sup>10</sup> Joe Jacob, *What Sites Were Used for Training Google Bard AI?*, MEDIUM (Feb. 11, 2023),  
[https://medium.com/@taureanjoe/what-sites-were-used-for-training-google-bard-ai-  
1216600f452d](https://medium.com/@taureanjoe/what-sites-were-used-for-training-google-bard-ai-1216600f452d).

<sup>11</sup> Ortiz, *supra* note 9.

<sup>12</sup> Jennifer Elias, *Google’s Newest A.I. Model Uses Nearly Five Times More Text Data for  
Training than Its Predecessor*, CNBC (May 17, 2023),  
[https://www.cnn.com/2023/05/16/googles-palm-2-uses-nearly-five-times-more-text-data-than-  
predecessor.html](https://www.cnn.com/2023/05/16/googles-palm-2-uses-nearly-five-times-more-text-data-than-predecessor.html).



1 but it is also built to generate a wide range of expression from shortform articles to books, chapters,  
2 mimicking expressive style, themes of the copyrighted works on which it was trained.

3 25. Research on LaMDA offers an overview of how the model’s 3.6 million tokens were  
4 sources to “achieve a more robust performance on dialog tasks:”

- 5 a. 50% dialog data from public forums;
- 6 b. 12.5% C4 data;
- 7 c. 12.5% code documents from sites related to programming...;
- 8 d. 12.5% Wikipedia (English);
- 9 e. 6.25% English web documents;
- 10 f. 6.25% Non-English web documents.<sup>13</sup>

11 26. Only items (b) “C4 data” and (e) “Wikipedia” are comprised of known data. The  
12 remaining 75% of the LaMDA dataset are ambiguous, generalized descriptors for websites and  
13 documents found across the internet. As one publication put it, “murky is the best word for  
14 describing the 75% of data that Google used for training LaMDA.”<sup>14</sup>

15 27. Of Google’s named sources, the C4 dataset contains copyrighted materials, including  
16 works found on pirating sites or nonconsenting digital libraries and subscription services. *See, infra*  
17 ¶¶ 39-48. While much of Bard/Gemini’s training material remains within a proverbial black box,  
18 the datasets that developers have made public reveal misappropriation of copyrighted works. *Id.*

19 28. Google released Gemini publicly on May 10, 2023, in over 180 countries and  
20 territories. Gemini quickly reached 142.6 million users the same month.<sup>15</sup> Google plans to expand  
21 to more countries, with an anticipated global reach of 1 billion users, or an eighth of all people  
22 worldwide.<sup>16</sup> Importantly, Google was determined to expedite the launch of its AI Products at the  
23 expense of creators’ exclusive rights—secretly harvesting millions of copyrighted materials from  
24 the internet without creators’ knowledge, consent, and consideration.

25 29. Google’s severe violation of artistic expression harms creators in a number of ways.  
26 Most clearly, they lose out on the monetary incentive that comes with being the holder of a

27 <sup>13</sup> Romal Thopplin, *et al.*, *LaMDA: Language Models for Dialog Applications*, GOOGLE (Feb. 10,  
28 2022), available at: <https://arxiv.org/pdf/2201.08239>

<sup>14</sup> Roger Montti, *Google Bard AI—What Sites Were Used To Train It*, SEARCH ENGINE JOURNAL  
(Feb. 10, 2023), <https://www.searchenginejournal.com/google-bard-training-data/478941/>

<sup>15</sup> *Id.*; David F. Carr, *As ChatGPT Growth Flattened in May, Google Bard Rose 187%*, SIMILARWEB:  
BLOG (June 5, 2023), <https://www.similarweb.com/blog/insights/ai-news/chatgpt-bard/>.

<sup>16</sup> Ritik Sharma, *23 Amazing Google Bard Statistics (Users, Facts)*, CONTENTDETECTOR.AI (June  
28, 2023), <https://contentdetector.ai/articles/google-bard-statistics>.

1 copyright. Google is not licensing the works as it is required, but instead is stealing them, by  
2 reproducing them and integrating them into their AI Products, without the consent of or  
3 compensation to the affected creators.

4 30. But ultimately, this disrespects the very notion of being an author, artist, or creator.  
5 These individuals devote a significant amount of time, energy, effort, creativity, and heart to develop  
6 their work. To illustrate, writing a book takes years and requires an author to jump through a variety  
7 of hurdles. They must complete a manuscript, find an agent, send their manuscript to an editor, and  
8 obtain a book deal with a publisher.<sup>17</sup> An author spends months to years to write a book. If an author  
9 were zealously writing 500 words a day for seven days a week, it would still take around five and a  
10 half months to complete a 300-page book.<sup>18</sup> It can then take anywhere from weeks to several years  
11 for an editor to complete their revisions of the book.<sup>19</sup> An author certainly does not undertake this  
12 extremely arduous process just for a massive corporation to take it for free.

13 31. When Google steals and reproduces the works without payment to the authors to make  
14 a multi-billion-dollar invention, it deprives authors of deserved royalties, undermining their  
15 financial stability and incentives to create new works.

16 32. Rather than licensing this valuable creative data from the owners, Defendant elected  
17 to systematically steal intellectual property of others, including personal and copyright information  
18 obtained without consent, and utilize datasets (such as the C-4 dataset) that are riddled with  
19 copyrighted and protected works, with the copyright symbol appearing more than 200 million times  
20 within the dataset.<sup>20</sup>

21 33. The law does not allow this kind of systematic infringement that Defendant has been  
22 and continues to commit.

23  
24 <sup>17</sup> *How Can I Get Published?* PENGUIN RANDOM HOUSE,  
25 <https://www.penguinrandomhouse.com/articles/how-can-i-get-published/> (last visited June 26,  
26 2024).

27 <sup>18</sup> *How Long Does It Take to Write a Book?* MASTERCLASS (Mar. 28, 2022),  
28 <https://www.masterclass.com/articles/how-long-does-it-take-to-write-a-book>.

<sup>19</sup> *How Can I Get Published?* *Supra* note 17.

<sup>20</sup> Kevin Schual, et. al., *Inside the secret list of websites that make AI like ChatGPT sound smart*,  
THE WASHINGTON POST (April 19, 2023),  
<https://www.washingtonpost.com/technology/interactive/2023/ai-chatbot-learning/>.

1           34. “Scraping involves the use of ‘bots,’ or robot applications deployed for automated  
2 tasks, which *scan and copy* the information on webpages then *store and index* the information.”<sup>21</sup>  
3 Scraping is stealing and infringing on the protected intellectual property of others. To do this,  
4 Google has to reproduce the works to then use them to build their AI Products. According to a  
5 computer science professor at the University of Oxford, the full extent of personal and copyrighted  
6 data taken by Defendant’s scraping is “unimaginable.”<sup>22</sup>

7           35. Authors, including Plaintiff, publish books with certain copyright management  
8 information. This information includes the book’s title, the ISBN number or copyright number, the  
9 author’s name, the copyright holder’s name, and terms and conditions of use. Most commonly, this  
10 information is found on the back of the book’s title page and is customarily included in all books,  
11 regardless of genre. The copyright symbol appears on the cover and initial pages of the books.  
12 Therefore, when Google downloads, copies, and otherwise reproduces creative works like  
13 Plaintiff’s book, to train its AI Products, it knows that the works are registered with the copyright  
14 office.

15           36. Google’s theft on a massive scale of copyrighted creative works without permission  
16 raises serious legal, moral, and ethical questions. Regulators and courts worldwide are seeking to  
17 crack down on AI companies “hoovering up content without consent or notice,”<sup>23</sup> but the response  
18 by Google and others has been to keep its training datasets largely secret. Google has not permitted  
19 any regulatory or other audit access.

20 ///

21 ///

22 ///

23 \_\_\_\_\_  
24  
25 <sup>21</sup> Brian Stuenkel, *Personal Information and Artificial Intelligence: Website Scraping and the*  
26 *California Consumer Privacy Act*, COLO. TECH. L. J. (Nov. 2, 2021),  
27 <https://ctlj.colorado.edu/?p=840>.

28 <sup>22</sup> Alex Hern & Dan Milmo, *I Didn’t Give Permission: Do AI’s Backers Care About Data Law Breaches?*, GUARDIAN (Apr. 10, 2023), <https://www.theguardian.com/technology/2023/apr/10/i-didnt-give-permission-do-ais-backers-care-about-data-law-breaches>.

<sup>23</sup> *Id.*

1                   **A. Google’s AI Product Development Depends on Stolen Vast Troves of**  
 2                   **Copyrighted Data, Including Plaintiff’s Book.**

3           37.     Around 50% of the C-4 dataset, created by Google in 2020, is taken from the Common  
 4     Crawl dataset.<sup>24</sup> The Common Crawl dataset is a massive collection of web pages and websites  
 5     consisting of petabytes of data collected over twelve (12) years, including raw web page data,  
 6     metadata extracts, and text extracts. The Common Crawl dataset is owned by a non-profit of the  
 7     same name, which has been indexing and storing as much of the internet as it can access, filing  
 8     away as many as 3 billion webpages every month, for over a decade.<sup>25</sup>

9           38.     The remaining 50 percent of the C-4 dataset is sourced from what Google vaguely  
 10    terms as “public forums.” The company has declined to clarify the specifics of what constitutes  
 11    these “public forums,” leaving users in the dark about the exact origins and nature of the data  
 12    influencing half of the AI’s training.<sup>26</sup>

13           39.     Importantly, the C-4 dataset is rife with copyrighted and protected works, with the  
 14    copyright symbol appearing more than 200 million times within the dataset.<sup>27</sup>

15           40.     In fact, the third largest site fueling the dataset is scribd.com, a subscription-based  
 16    digital library with sixty (60) million e-books and audio books—that compensates authors using a  
 17    revenue sharing model based on the number of reads their work gets.<sup>28</sup> Word for word excerpts  
 18    from Plaintiff Leovy’s copyrighted work appears on scribd.com. There is no indication Scribd  
 19    consented to Google’s mass misappropriation of copyrighted works on this website, and certainly  
 20    Plaintiff Leovy nor any other author consented to Google’s use of this material, nor were they  
 21    compensated. Rather, Google has engaged in unauthorized accessing and theft of copyrighted and  
 22    restricted materials.

23           \_\_\_\_\_  
 24    <sup>24</sup> *Id.*; Katyanna Quach, *4chan and Other Web Sewers Scraped Up Into Google’s Mega-Library*  
 for Training ML, THE REGISTER (Apr. 20, 2023),

[https://www.theregister.com/2023/04/20/google\\_c4\\_data\\_nasty\\_sources/](https://www.theregister.com/2023/04/20/google_c4_data_nasty_sources/).

25    <sup>25</sup> James Bridle, *The Stupidity of AI*, GUARDIAN (Mar. 16, 2023),

<https://www.theguardian.com/technology/2023/mar/16/the-stupidity-of-ai-artificial-intelligence-dall-e-chatgpt>.

26    <sup>26</sup> Roger Montti, *Google Bard AI: What Sites Were Used to Train It*, SEARCH ENGINE J. (Feb. 10,  
 27    2023), <https://www.searchenginejournal.com/google-bard-training-data/478941/>.

27    <sup>27</sup> Schaul, *supra* note 20.

28    <sup>28</sup> *Id.*; Omar, *Scribd Review: Scribd Membership Options, Pros, Cons, and Pricing*, OJ DIGIT.  
 SOLUTIONS (Last updated April 29, 2023), <https://ojdigitalsolutions.com/scribd-review/>.

1           41. Worse, the dataset Google used also contains data from “b-ok.org” a “notorious  
2 market for pirated e-books,” as well as “[a]t least 27 other sites identified by the U.S. government  
3 as markets for piracy and counterfeits.”<sup>29</sup>

4           42. “B-ok.org,” also known as “Z Library,” is “[t]he world’s largest [illegal] ebook library  
5 and digital library.”<sup>30</sup> It is a pirated cite that is being prosecuted for criminal copyright infringement,  
6 and the books that illegally appeared on Z Library, like Plaintiff Leovy’s book in its entirety, were  
7 not legally authorized to be available for distribution or copying.

8           43. Google nonetheless took Plaintiff Leovy’s book in its entirety, unlawfully  
9 misappropriated and reproduced it and utilized it to train Google’s AI Products, as it did with  
10 countless other author’s copyrighted works taken from this and myriad other websites.

11           44. Additionally, there is a trove of personal blogs represented in the misappropriated  
12 data—more than half a million, including the tens of thousands of blogs hosted on Medium, a  
13 website especially popular with authors and other content creators. Blogs written on WordPress,  
14 Tumbler, Blogspot and Live Journal were also among the materials misappropriated by Google.

15           45. Google also misappropriated copyrighted information from popular crowdfunding  
16 and creative websites, Kickstarter and Patreon, giving Gemini access to thousands of artists’ and  
17 creators’ proprietary marketing materials, “raising concerns [Gemini] may copy this work in  
18 suggestions to users.”<sup>31</sup>

19           46. Websites and platforms like Patreon are dedicated to getting creators paid for their  
20 exclusive content.<sup>32</sup> Google ignored this and stole creators’ exclusive works anyway.

21           **B. Google’s Revised Privacy Policy Purports to Give it “Permission” to Take**  
22           **Anything Shared Online to Train and Improve Its AI Products, Including**  
23           **Copyrighted Information.**

24           47. On July 1, 2023, Google quietly amended its privacy policy to openly assert that it  
25 “collects” publicly available information from the web to train its AI Products, including “Gemini”

26 <sup>29</sup> Schaul, *supra* note 20.

27 <sup>30</sup> Beinginstructor, *ZLibrary — The world’s largest ebook library*, MEDIUM (Feb. 16, 2023),  
28 <https://medium.com/@beinginstructor/zlibrary-the-worlds-largest-ebook-library-dfb933762cfc>.

<sup>31</sup> Schaul, *supra* note 20.

<sup>32</sup> *The Story of Patreon*, PATREON, <https://www.patreon.com/about> (last visited June 27, 2024).

1 and “Cloud AI.”<sup>33</sup> Given that Google had been doing this theft on a massive scale in secret for years,  
 2 this disclosure was long overdue. But it was also alarming because it solidified as corporate “policy”  
 3 Google’s disregard for the internet community as a whole, but particularly for authors and creators,  
 4 and reflected its intent to continue exploiting for commercial gain copyrighted works that are  
 5 “publicly available online.”

6 Figure 3

7 **publicly accessible sources**

8  
 9 For example, we may collect information that’s publicly available online or from other  
 10 public sources to help train Google’s language AI models and build products and features  
 11 like Google Translate, Bard, and Cloud AI capabilities. Or, if your business’s information  
 12 appears on a website, we may index and display it on Google services.

13 48. The idea that Google believes all publicly available creative works on the internet are  
 14 fair game for it to take, commercially misappropriate, and build AI Products has shocked and  
 15 angered the creators. As one bestselling author, Alexander Chee, explained, “There’s no urgent need  
 16 to AI to write a novel. The only people who might need that are the people who object to paying  
 17 writers what they’re worth.”<sup>34</sup> The CEO of the Author’s Guild plainly stated, “It’s not fair to use  
 18 our stuff in your AI without permission or payment. So please start compensating and talking to  
 19 us.”<sup>35</sup>

20 49. Responding to the backlash, Google announced it will host a public forum to discuss  
 21 what data collection and protection practices should look like in the new AI era.<sup>36</sup> But as many  
 22 internet users noted, it is a little too late for that now that Google has already taken and  
 23

24 <sup>33</sup> Jess Weatherbed, *Google Confirms it’s Training Bard on Scraped Web Data, Too*, THE VERGE  
 25 (July 5, 2023), <https://www.theverge.com/2023/7/5/23784257/google-ai-bard-privacy-policy-train-web-scraping>.

26 <sup>34</sup> Chloe Veltman, *Thousands of Authors Urge AI Companies to Stop Using Work Without  
 27 Permission*, NPR (July 17, 2023), <https://www.npr.org/2023/07/17/1187523435/thousands-of-authors-urge-ai-companies-to-stop-using-work-without-permission>.

28 <sup>35</sup> *Id.*

<sup>36</sup> Matt G. Southern, *Google Calls for Public Discussion on AI Use of Web Content*, SEARCH  
 ENGINE J. (July 7, 2023), <https://www.searchenginejournal.com/google-calls-for-public-discussion-on-ai-use-of-web-content/491053/>.



1 misappropriated copyrighted works. In the words of one, Google is essentially saying to the world:  
2 “Now that we’ve already trained our LLMs on all your proprietary and copyrighted content, we will  
3 finally start thinking about giving you a way to opt out of any of your future content being used to  
4 make us rich.”<sup>37</sup> Google’s willingness to *potentially* allow authors to “opt out” of the future use  
5 content does not change the fact that Google has illegally copied and misappropriated the  
6 copyrighted works, and continues to engage in copyright infringement when it engages in further  
7 theft of the same databases and websites that contain pirated works.

8 50. Defendant’s illegal and invasive data misappropriation and infringement practices  
9 have also led social platforms that contain copyright protected content, like Twitter and Reddit, to  
10 enact more stringent measures in an effort to protect the rights and data of its millions of users.<sup>38</sup>  
11 But these anti-scraping modifications stand to negatively impact use of the internet for everyone.  
12 For example, now the public cannot view tweets unless they are logged in to Twitter and are limited  
13 in how many tweets they can view in one day.

14 51. These negative impacts to the internet at large underscore the unfortunate ripple  
15 effects of Google’s misconduct.<sup>39</sup> Unless Google and other AI giants like it are ordered to stop the  
16 illegal theft of copyrighted material, other websites might be forced to similarly limit access to the  
17 public.

18 52. As one commentator observed, “should sites really have to wall off their mountains  
19 of text so that AI companies can’t gobble it up and use it to build AI? That makes no sense.”<sup>40</sup> If  
20 this were to happen at scale, it would forever change how the internet works, limiting its utility for  
21 millions of good faith users who do not want to steal data, but simply engage with it legally in  
22 accordance with a site’s terms of use and the property interests of the content creators themselves.

23 53. Moreover, the new policy does not except use of copyrighted (or any other) material

24 \_\_\_\_\_  
25 <sup>37</sup> *Id.*

26 <sup>38</sup> *Musk Says Twitter Will Limit How Many Tweets Users Can Read*, REUTERS (July 1, 2023),  
<https://www.reuters.com/technology/musk-says-twitter-applies-temporary-limit-address-data-scraping-system-2023-07-01/>.

27 <sup>39</sup> Cory Woodroof, *Twitter Users Were Furious After the Website Temporarily Applied a Reading  
Limit*, USA TODAY (July 1, 2023), <https://ftw.usatoday.com/lists/twitter-rate-limit-exceeded-elon-musk-angry-reactions>.

28 <sup>40</sup> Josh Marshall, *Twitter, Musk and the Great AI Land Grab*, TALKING POINTS MEMO (July 6,  
2023), <https://talkingpointsmemo.com/edblog/twitter-musk-and-the-great-ai-land-grab>.



1 from being included in its scraped data pool further exposing Google’s disregard for intellectual and  
2 other property rights while also undermining the policies of various publicly accessible websites,  
3 which explicitly prohibit *any* data collection or web scraping for the purpose of training AI models.

4 54. Now that Google has essentially claimed ownership rights over anything online, there  
5 is reason to believe that Google will violate the copyright interests of millions more. Indeed, a  
6 massive portion of Defendant’s misappropriation to date already includes the unauthorized and  
7 widespread theft and copying of copyrighted works extending across a wide spectrum of industries  
8 that depend on creative and unique content creation.

9 55. Instead of competing fairly, Defendant illegally copied the unique works of millions  
10 of creators to develop and “train” its AI technology, without consent, credit, or fair compensation.  
11 This unauthorized theft and usage of copyrighted content stands in stark violation of creators’  
12 exclusive rights under copyright law.

13 56. Considering the magnitude and scale of the copyright violations to date, along with  
14 the likelihood that these violations will continue to increase exponentially, content creators will be  
15 dissuaded from investing in the considerable costs of producing unique content in electronic  
16 formats. This not only threatens to drastically reshape online accessibility of paid, restricted  
17 materials, but also imposes economic harm on a substantial number of content creators.

18 57. Despite the existence of numerous lawful ways to acquire training data, Defendant  
19 purposely elected to bypass the legal routes, opting instead to pillage the internet for copyrighted  
20 works. The resulting impact has not only infringed upon the rights of countless creators but has  
21 created an environment that ultimately discourages creativity and innovation.

22 58. The AI Products also dramatically undercuts the commercial market for books and  
23 other works already created by radically altering the perceived incentives for anyone to purchase  
24 the stolen works going forward. This further harms millions of authors and creators in the form of  
25 lost profits and otherwise.

26 **C. Google Uses This Stolen Data to Profit by the Billions.**

27 59. Google’s unlawful theft of copyrighted material, at no cost to train its AI technology,  
28 has and will continue to unjustly enrich Google. For example, Google announced Gemini, which at

1 the time was called “Bard,” on February 6, 2023, and the very next day Alphabet Inc.’s market  
2 capitalization increased to 1.37 trillion, reaching 1.62 trillion in June of 2023—its highest market  
3 capitalization in the past year.<sup>41</sup>

4 60. Only a few months after announcing Gemini and in the wake of the AI frenzy, Google  
5 co-founders Larry Page and Sergey Brin experienced a combined wealth increase of over \$18 billion  
6 as the company revealed a revamped AI powered search engine.<sup>42</sup> Page’s net worth increased by  
7 \$9.4 billion to \$106.9 billion, while Brin’s increased by \$8.9 billion to \$102.1 billion.<sup>43</sup>

8 61. This is far from a short-lived AI inspired spike. Google cleverly monetizes its AI  
9 Products and fails to meaningfully disclose that Google uses valuable copyrighted material to  
10 enhance other Google products and services *and net billions*.<sup>44</sup> Gemini sweeps in profit for Google  
11 on the backs of copyrighted authors and creators without paying any licensing fees for the  
12 unauthorized reproduction of their works to train Gemini.

13 62. Google AI’s DeepMind is alone now worth around \$55 million,<sup>45</sup> yet the individuals  
14 and companies that produced the copyrighted material Google illegally copied and misappropriated  
15 from the internet have not been compensated. This Action seeks to change that, and in the process,  
16 protect the property and privacy rights of millions.

#### 17 **D. Creators are Outraged by Google’s Theft-Based Training Model**

18 63. Google has continued to harvest mass amounts of copyrighted material despite an  
19 outpour of public outrage. Specifically, creators have recognized and expressed discontent with  
20 Google’s problematic business model, which allows it to unfairly profit off artists, authors, and  
21 content creators, and that forces everyone, whether they want to or not, to contribute to building

22 \_\_\_\_\_  
23 <sup>41</sup> *Google Announces Bard, Its Rival to Microsoft-Backed ChatGPT*, FORBES (Feb. 8, 2023),  
24 <https://www.forbes.com/sites/qai/2023/02/08/google-announces-bard-its-rival-to-microsoft-backed-chatgpt/?sh=29ed0fd93791>; *Alphabet Market Cap 2010-2023*, MACROTRENDS,  
25 <https://www.macrotrends.net/stocks/charts/GOOGL/alphabet/market-cap> (last visited June 27,  
26 2024).

<sup>42</sup> Biz Carson, *Google Co-Founders Gain \$18 Billion as AI Boost Lifts Stock*, BLOOMBERG (May  
12, 2023), <https://www.bloomberg.com/news/articles/2023-05-12/google-co-founders-gain-17-billion-as-ai-boost-lifts-stock>.

<sup>43</sup> *Id.*

<sup>44</sup> *Bard Privacy Notice*, BARD, <https://support.google.com/bard/answer/13594961?hl=en> (last  
27 updated May 29, 2024).

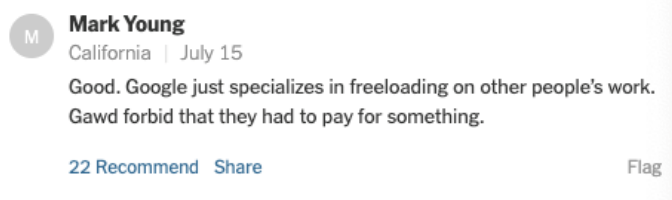
<sup>45</sup> *DeepMind Net Worth*, PEOPLE AI, <https://peopleai.com/fame/identities/deepmind> (last visited  
28 Jan. 1, 2024).

1 untested and volatile technology that violates copyright laws.

2 64. Users are rightfully upset that the content they invest their time and energy into, and  
3 in all cases, which is intended for specific audiences and purposes is being used to create a  
4 multibillion-dollar franchise that they will never see a dime of. One X user shared, “Authors – your  
5 creative work is valuable. It deserves protection. You have the right to control what happens to it.  
6 Google is allegedly data scraping all the documents in google docs to train their AI. This includes  
7 your work! #writingcommunity.”<sup>46</sup>



8  
9  
10  
11  
12  
13  
14  
15 65. One New York Times reader expressed a similar sentiment: “Google just specializes  
16 in freeloading on other people’s work. Gawd forbid they had to pay for something.”<sup>47</sup>



17  
18  
19  
20  
21  
22 66. Similarly, another New York Times reader added: “Once again, capitalism proves it’s  
23 obsessed with the idea of a zero-expense operation – if it can get what it wants for free and only  
24

25 <sup>46</sup> Kelsey Brownlee (@\_kelseybrownlee), X (July 14, 2023),  
26 [https://x.com/\\_kelseybrownlee/status/1679954300376686594?s=46&t=HHkRbC2AV14Ias31BERw9g](https://x.com/_kelseybrownlee/status/1679954300376686594?s=46&t=HHkRbC2AV14Ias31BERw9g).

27 <sup>47</sup> Sheera Frenkel & Stuart A. Thompson, ‘Not for Machines to Harvest’: Data Revolts Break Out  
28 *Against A.I.*, THE N. Y. TIMES, (July 15, 2023)  
<https://www.nytimes.com/2023/07/15/technology/artificial-intelligence-models-chat-data.html#commentsContainer>. Commenter: Mark Young.

1 collect revenues from customers, that is what it could consider nirvana. The prospect of assuming  
 2 anything publicly visible to be free of charge, and then cutting creators out of any receipts, is what  
 3 especially has creators rightfully up in arms.”<sup>48</sup> The reader bluntly added, “You know who else  
 4 collects money without giving anything back in return? Robbers.”<sup>49</sup>



**IlliniWatcher**

Houston | July 15

I've been saying it since the start of the AI hype - the entire industrial world is about to get an important lesson on ethics. And I've worked in the IT industry for decades, so I'm a bit closer to the action than those who get their info on tech from Hollywood and streaming series.

Once again, capitalism proves it's obsessed with the idea of zero expense operation - if it can get what it wants for free and only collect revenues from customers, that is what it would consider nirvana. The prospect of assuming anything publicly visible to be free of charge, and then cutting creators out of any receipts, is what especially has creators rightfully up in arms.

You know who else collects money without giving anything back in return? Robbers. Robbers only take, expecting they won't get caught, and pocket whatever they can get from the unsuspecting.

A lot of business models MUST change. The suits at the top have obscene compensation packages while the vast majority of the rank and file - the talent - gets edged out of the picture. It's also happening in entertainment (writers and, as of this past week, actors), shipping (witness the UPS brouhaha) and retail coffee (exhibit A: Starbucks).

All it comes down to is learning to share the wealth - and the respect - with talent and its many creators.

[34 Recommend](#) [Share](#)

[Flag](#)

22 67. Another reader shared a digestible analogy that proves that users can see through  
 23 Google’s mystique. “But if I said ‘here is the work I created in the style of JK Rowling!’ and it was  
 24 just mashed together and reworded sentences from the Harry Potter books, I’d be laughed out of the  
 25 room.”<sup>50</sup> Despite AI’s smoke-and-mirrors, users can see that big tech’s technological advancement  
 26 is nothing more than wide-scale theft.

27 <sup>48</sup> *Id.* Commenter: IlliniWatcher.

28 <sup>49</sup> *Id.*

<sup>50</sup> *Id.* Commenter: Cody.

**Cody**

British Columbia | July 15

People seriously need to think through on their own whether they actually believe what AI is doing is impressive or cool or helpful; so many people are just repeating what they've heard others say and calling the technology "powerful" and "impressive" out of fear of being labelled a luddite or out of touch. News outlets are breathlessly doing free advertising for these companies by talking about their "impressive" capabilities.

But if I said "here is the work I created in the style of JK Rowling!" and it was just mashed together and reworded sentences from the Harry Potter books, I'd be laughed out of the room. But for some reason people think its incredible when the chatbot does it.

Oh but it's just in its infancy and it will create truly impressive works of literature one day right? Get back to me when it does. For 20 years people have been saying self-driving cars and trucks will put delivery drivers and truckers out of work, and all I see are news articles about trucker shortages.

68. Artists, creators, and writers have voiced that they feel particularly threatened by Defendant's data-theft tactics. Many of these users' livelihoods are dependent on sharing their content on the internet. When they discovered that creations that they poured their expertise into were being stolen, illegally copied, and used to train AI products—without any form of acknowledgement or compensation—they were rightfully upset.

69. In fact, The Author's Guild shared an open letter they wrote to AI companies.<sup>51</sup> The letter begged that these companies, as the "leaders of AI" take steps to "mitigate the damage to [their] profession" caused by data scraping and AI training.<sup>52</sup> Collectively, the authors asked that AI companies, including Google, "Compensate writers fairly for the past and ongoing use of our works in your generative AI programs."<sup>53</sup>

70. Eva Toorenent, an illustrator who serves as the Netherland's advisor for the European Guild for Artificial Intelligence, argued that "[AI models] have sucked the creative juices of millions

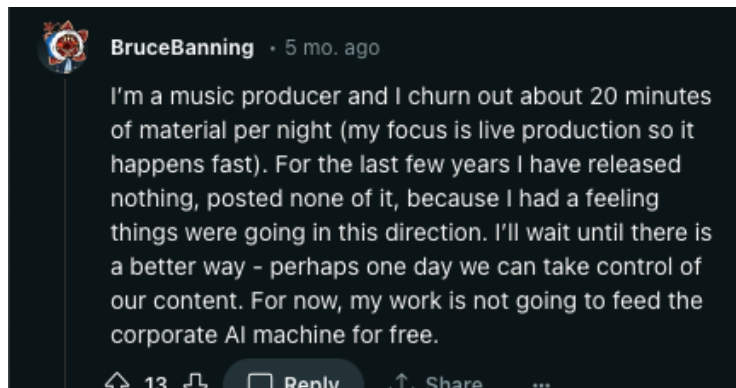
<sup>51</sup> The Author's Guild, *Open Letter to Generative AI Leaders*, <https://actionnetwork.org/petitions/authors-guild-open-letter-to-generative-ai-leaders> (last visited June 27, 2024).

<sup>52</sup> *Id.*

<sup>53</sup> *Id.*

1 of artists.”<sup>54</sup> Molly Crabapple, a writer and artist, similarly shared, “To see corporations scrape our  
2 style and then attempt to replace us with bastardized versions of our own work is beyond  
3 disgusting.”<sup>55</sup>

4 71. The threat of AI companies, like Defendant’s, misappropriating and reproducing  
5 users’ content has caused some creators to refrain from posting their content altogether. One Reddit  
6 user shared, “For the last few years I have released nothing,” referring to the music he produces.<sup>56</sup>  
7 He added, “perhaps one day we can take control of our content. For now, my work is not going to  
8 feed the corporate AI machine for free.”<sup>57</sup>



15 72. Absent injunctive relief sought herein, Plaintiff and the Class will continue to not  
16 freely contribute online as they might for fear of losing control of their data.

17 **E. The Public is Outraged by the Lack of Respect for Autonomy in the**  
18 **Copyright Space, and AI Developments Writ Large**

19 73. The US Copyright Office opened a public comment period on August 30, 2023,  
20 concerning the use of copyrighted data to train AI models, including the violation of publicity  
21 rights.<sup>58</sup>

22 74. Several individuals noted the glaring invasion of privacy that AI companies are

23

24 <sup>54</sup> Kate Knibbs, *A new Tool Helps Artists Thwart AI—With a Middle Finger*, WIRED (Oct. 12, 2023), <https://www.wired.com/story/kudurru-ai-scraping-block-poisoning-spawning/>.

25 <sup>55</sup> *Id.*

26 <sup>56</sup> Bruce Banning, *Google's policy update confirms that all your posted content will be utilized for AI training*, REDDIT, (June 2023), [https://www.reddit.com/r/technews/comments/14qe9tm/googles\\_policy\\_update\\_confirms\\_that\\_all\\_your/?sort=top](https://www.reddit.com/r/technews/comments/14qe9tm/googles_policy_update_confirms_that_all_your/?sort=top).

27 <sup>57</sup> *Id.*

28 <sup>58</sup> Emilia David, *US Copyright Office Wants to Hear What People Think About AI and Copyright*, THE VERGE (Aug. 29, 2023), <https://www.theverge.com/2023/8/29/23851126/us-copyright-office-ai-public-comments>.



1 engaging in, beyond just copyright. For example, one commenter wrote: “The current practice of  
2 using AI to create art/text/video/etc by feeding it people’s personal information, conversations, and  
3 artistic work seems like both **obvious plagiarism/copyright infringement**, and a major breach of  
4 privacy for every person living in this country.”<sup>59</sup>

5 75. Another commenter shared, “**Never have I consented to have any of the work I’ve**  
6 **posted online be used to fuel an AI engine, and I certainly don’t consent to allowing the people**  
7 **behind said AI and scrapping to profit off of my work or other things I’ve posted.** I do not feel  
8 comfortable having personal work used to power an engine made to generate profit, of which I will  
9 never see a penny of . . . It’s violating our trust and privacy, not to mention the amount of  
10 copyrighted works it has scraped from online pdfs and others sources to build this AI. **This isn’t**  
11 **legal, as it’s directly stealing and profiting off of stolen content, not adding anything new to**  
12 **it.**”<sup>60</sup>

13 76. The comments exhibited an overwhelming level of infuriation over the sad reality that  
14 the creative works of millions are being exploited:

15 “As a working professional artist, where my entire income rests upon  
16 my artwork, I feel like it is not okay for generative ai companies to be  
17 disguising themselves as nonprofit and data laundering my artwork for  
18 their profit. I would never opt in to companies like this even if I were to  
19 be compensated fairly. I do not want my artwork to be trained for Ai. I  
20 do not want any of my personal information to be training any sort of  
21 data set. My job is literally be replaced right now as we speak because  
22 everyone is ‘having fun’ at the expense of my livelihood. Please do not  
23 continue letting this companies slide.”<sup>61</sup>

## 24 **F. Online News and Media Businesses are Taking Action Against Google’s** 25 **Unauthorized Infringement**

26 77. Much like the average internet user, many online news and media websites are  
27 concerned that Defendant is stealing data to train their AI models.

28 <sup>59</sup> *Comment from Clorite, Katelyn*, U.S. COPYRIGHT OFFICE (Oct. 30, 2023),  
<https://www.regulations.gov/comment/COLC-2023-0006-1003> (emphasis added).

<sup>60</sup> *Comment from Anonymous*, U.S. COPYRIGHT OFFICE (Oct. 31, 2023),  
<https://www.regulations.gov/comment/COLC-2023-0006-5235>.

<sup>61</sup> *Comment from Chan, Maggie*, U.S. COPYRIGHT OFFICE (Oct. 30, 2023),  
<https://www.regulations.gov/comment/COLC-2023-0006-0347>.



1           78. To combat unlicensed data collection, hundreds of publishers are trying to block AI  
2 web-crawlers from scanning their websites. Included in the list of media giants that have inserted  
3 code in an attempt to block web crawlers, on a go forward basis, are the New York Times, CNN,  
4 Reuters, Disney, Bloomberg, The Washington Post, ABC News, ESPN, and Insider.

5           79. There is increasing concern that generative AI, if it continues to grow at this rate,  
6 could greatly impact the publishing industry and even go as far as to put some newsrooms out of  
7 business. This would be ironic, given that AI's growth is and has been dependent on stealing  
8 information from these very sources.

9           80. News stories are a critical resource in developing generative AI. These companies'  
10 outrage demonstrates that they recognize the value of their content and believe that they should not  
11 be allowing AI web-crawlers to capitalize on that their content without paying for it in the first  
12 place. Similar to the reactions of average internet users, these companies' response demonstrates  
13 the overarching anger towards Defendant's unfair and anticompetitive practices—spanning across  
14 the entire internet food-chain.

## 15           **II. DEFENDANT'S CONDUCT VIOLATES ESTABLISHED COPYRIGHT LAWS.**

16           81. Defendant's theft and misappropriation violated the copyright rights of all individuals  
17 whose creative content was copied and then incorporated into Defendant's AI Products.

18           82. Defendant's illegal stealing and reproducing of copyrighted works was done largely  
19 in secret, without consent from or consideration to any creator whose copyrighted material was  
20 reproduced.

21           83. Creators retain copyright interests over their unique and original content posted (or at  
22 times pirated and illegally displayed) online. This content includes text, images, music, video  
23 content, and other forms of creative expression, all of which fall under the purview of copyright  
24 law.

25           84. Defendant's unauthorized theft, reproduction, and use of these copyrighted materials  
26 constitutes infringement because Defendant copies and downloads the intellectual property to then  
27 use it to build and train its AI Products. As an illustrative example, the unauthorized collection and  
28 use of copyrighted literary works in training Gemini not only infringes on the rights of the producers

1 but also damages the intrinsic value of the copyrighted works.

2 85. Copyright protection incentivizes creativity and original content creation. Copyright  
3 holders have exclusive rights to reproduce their work in different formats, commercially exploit it,  
4 create derivative works, and display or perform the work publicly. Thus, when copyrighted work is  
5 co-opted without permission or compensation, as in the case of Defendant's massive theft and  
6 infringement, it severely undermines the fundamental principles of copyright law.

7 86. Further, the practice of illegal theft and infringement of works through web scraping  
8 effectively nullifies the concept of "fair use," a critical aspect of copyright law designed to allow  
9 limited use of copyrighted material without permission for purposes like commentary, criticism,  
10 news reporting, and scholarly reports. *See McGucken v. Pub Ocean Limited*, 42 F.4th 1149 (9th Cir.  
11 2022). Defendant's wholesale collection and use of copyrighted material, with no option for  
12 copyright owners to opt out, far exceeds any reasonable interpretation of "fair use." *See VHT v.*  
13 *Zillow Group*, 918 F.3d 723, 743 (9th Cir. 2019); *accord Worldwide Church of God v. Phila. Church*  
14 *of God, Inc.*, 227 F.3d 110, 1118 (9th Cir. 2000) ("[C]opying an entire work militates against a  
15 finding of fair use.").

16 87. The non-consensual aggregation and usage of copyrighted materials disrupts the  
17 balance between content creators and consumers that copyright law intends to foster. When original  
18 content is unfairly utilized in this manner, it discourages creators from investing time, effort, and  
19 resources into creating new content.

20 88. By using such works as training fodder for its AI, Defendant is not just using these  
21 works in an unauthorized manner, but also illegally profiting from infringement. Plaintiff and Class  
22 Members have not consented to such exploitation of their copyrighted works. It is only through  
23 legal action that the rights of content creators can be protected, and their original works safeguarded  
24 against such egregious misuse.

25 89. While the past, and ongoing, misappropriation of valuable copyrighted material is bad  
26 enough, AI Products like Gemini also stand to altogether eliminate future income for millions, due  
27 to the widespread unemployment AI and loss of value for intellectual property it expected to cause  
28 over time. No one has consented to the use of their copyrighted materials in a manner that not only

1 violates copyright laws but that also may build this destabilized future of social unrest and  
2 worsening poverty for everyday people, while the pockets of Google are lined with profit.

### 3 CLASS ALLEGATIONS

4 90. **Class Definition:** Plaintiff brings this action pursuant to Federal Rules of Civil  
5 Procedure Sections 23(b)(2), 23(b)(3), and 23(c)(4), on behalf of Plaintiff and the Class defined as  
6 follows:

7 **Copyright Class:** All persons in the United States who own a United  
8 States copyright in any work that was used as training data for  
9 Defendant's Products.

10 91. **The following people are excluded from the Class:** (1) any Judge or Magistrate  
11 presiding over this action and members of their judicial staff and immediate families; (2) Defendant,  
12 Defendant's subsidiaries, parents, successors, predecessors, and any entity in which the Defendant  
13 or its parents have a controlling interest and its current or former officers and directors; (3) persons  
14 who properly execute and file a timely request for exclusion from the Class; (4) persons whose  
15 claims in this matter have been finally adjudicated on the merits or otherwise released; (5) Plaintiff's  
16 counsel and Defendant's counsel; and (6) the legal representatives, successors, and assigns of any  
17 such excluded persons. Furthermore, the Class excludes any works which currently are in public  
18 domain.

19 92. Plaintiff reserves the right under Federal Rule of Civil Procedure 23 to amend or  
20 modify the Class to include a broader scope, greater specificity, further division into subclasses, or  
21 limitations to particular issues. Plaintiff reserves the right under Federal Rule of Civil Procedure  
22 23(c)(4) to seek certification of particular issues.

23 93. The requirements of Federal Rules of Civil Procedure 23(a), 23(b)(2), and 23(b)(3)  
24 are met in this case.

25 94. The Fed. R. Civ. P. 23(a) elements of Numerosity, Commonality, Typicality, and  
26 Adequacy are all satisfied.

27 95. **Ascertainability:** Membership of the Class is defined based on objective criteria, and  
28 individual members will be identifiable from Defendant's records, records of other Google  
products/services, self-identification methods, or other means. Defendant's records are likely to

1 include massive data storage, user accounts, and data gathered directly from the affected members  
2 of Class.

3 96. **Numerosity:** The precise number of Class Members is not available to Plaintiff, but  
4 it is clear that individual joinder is impracticable. Millions of copyright holders have been victims  
5 of Defendant's unlawful and unauthorized infringement and theft on massive scale. Class Members  
6 can be identified through Defendant's records of copyrighted works, records from the U.S.  
7 Copyright Office, or by other means, including but not limited to self-identification.

8 97. **Commonality:** Commonality requires that the Class Members allege claims which  
9 share common contention such that determination of its truth or falsity will resolve an issue that is  
10 central to the validity of each claim in one stroke. Here, there is a common contention for all Class  
11 Members are as follows:

- 12 a) Whether Leovy owned copyright in her book that was scraped, copied, and used to  
13 train Defendant's AI Products.
- 14 b) Whether Defendant's conduct constitutes an infringement of the copyrights held by  
15 Plaintiff Leovy and the Class in their respective works;
- 16 c) Whether Defendant's reproduction of the copyrighted works constitutes a copyright  
17 infringement;
- 18 d) Whether Defendant's copying and/or reproduction of the copyrighted works  
19 constitutes fair use;
- 20 e) Whether Defendant's violation of Class' and Plaintiff's exclusive rights under  
21 copyright law entitles them to damages, including statutory damages, and the  
22 amount of statutory damages;
- 23 f) Whether Defendant acted willfully with respect to the copyright infringements;

24 98. **Typicality:** Plaintiff's claims are typical of the claims of other Class Members in that  
25 Plaintiff and the Class Members sustained damages arising out of Defendant's uniform wrongful  
26 conduct and data collecting practices, and use of such data in an attempt to train the AI Products,  
27 and further develop the Products.

28 ///

1           **99. Adequate Representation:** Plaintiff will fairly and adequately represent and protect  
2 the interests of the Class Members. Plaintiff's claims are made in a representative capacity on behalf  
3 of the Class Members. Plaintiff has no interests antagonistic to the interests of the other Members  
4 of Class. Plaintiff has retained competent counsel to prosecute the case on behalf of herself and the  
5 Class. Plaintiff and Plaintiff's counsel are committed to vigorously prosecuting this action on behalf  
6 of the Class Members.

7           100. The declaratory and injunctive relief sought in this case includes, by way of example  
8 and without limitation:

- 9           a) Implementation of Accountability Protocols that hold Defendant responsible for  
10 Products' actions by barring any use of the protected materials until Plaintiff and  
11 Class Members are fairly compensated for the stolen and copied protected  
12 materials, and are compensated for the ongoing use for their intellectual property;
- 13           b) Implementation of Accountability Protocols that hold Defendant responsible for  
14 ensuring that during any web scraping it instructs web crawlers to avoid (a)  
15 websites/datasets that are known for containing pirated materials (such as Z-  
16 library and datasets containing Z-Library); (b) requiring that it reviews copyright  
17 notices and terms of the websites/databases which it uses for training to ensure  
18 that copyright protected materials are not within the training data sets; (c) limit  
19 scraping only to websites/datasets which contain materials within the public  
20 domain; (d) limit web scraping to datasets/websites for which Defendant has  
21 provided an accepted valuable consideration to authors and creator;
- 22           c) Requiring Defendant to allow Product users and everyday internet users to opt out  
23 of all collection/copying of their protected works, and stop the illegal taking of  
24 protected works, delete any ill-gotten protected materials, or the algorithms which  
25 were built on the stolen data;
- 26           d) Confirmation that Defendant has deleted, destroyed, and purged the copyrighted  
27 materials of all relevant class members unless Defendant can confirm that it has  
28 properly licensed such materials; and

- 1 e) Requiring all further and just corrective action, consistent with permissible law  
2 and pursuant to only those causes of action so permitted.

3 101. **This case also satisfies Fed. R. Civ. P. 23(b)(3) - Predominance:** There are many  
4 questions of law and fact common to the claims of Plaintiff and Class Members, and those questions  
5 predominate over any questions that may affect individual Class Members. Common questions  
6 and/or issues for Class members include the questions listed above in *Commonality*, and also  
7 include, but are not necessarily limited to the following:

- 8 a) Whether Defendant violated Plaintiff's and Class Members' exclusive rights under  
9 copyright laws;
- 10 b) Whether Plaintiff and Class members are entitled to actual damages, enhanced  
11 damages, statutory damages, restitution, disgorgement, and other monetary  
12 remedies provided by equity and law;
- 13 c) Whether injunctive and declaratory relief and other equitable relief is warranted.

14 102. **Superiority:** This case is also appropriate for class certification because class  
15 proceedings are superior to all other available methods for the fair and efficient adjudication of this  
16 controversy, as joinder of all parties is impracticable. The damages suffered by individual Class  
17 Members will likely be relatively small, especially given the burden and expense of individual  
18 prosecution of the complex litigation necessitated by Defendant's actions. Thus, it would be  
19 virtually impossible for the individual Class Members to obtain effective relief from Defendant's  
20 misconduct. Even if Class Members could mount such individual litigation, it would still not be  
21 preferable to a class action, because individual litigation would increase the delay and expense to  
22 all parties due to the complex legal and factual controversies presented in this Complaint. By  
23 contrast, a class action presents far fewer management difficulties and provides the benefits of single  
24 adjudication, economy of scale, and comprehensive supervision by a single Court. Economies of  
25 time, effort, and expense will be enhanced, and uniformity of decisions ensured.

26 103. Likewise, particular issues under Rule 23(c)(4) are appropriate for certification  
27 because such claims present only particular, common issues, the resolution of which would advance  
28 the disposition of this matter and the parties' interests therein.

1 **COUNT ONE**

2 **DIRECT COPYRIGHT INFRINGEMENT**

3 104. Plaintiff Leovy, individually and on behalf of the Class, herein repeats, realleges, and  
4 fully incorporates all allegations in all preceding paragraphs.

5 105. Copyrights are the legal title to intellectual property by which creators of original  
6 works (such as books, photographs, videos etc.) protect their moral, economic, and legal rights. The  
7 importance of copyrighted works is enshrined in the U.S. Constitution, which expressly gave  
8 Congress the power to “promote the Progress of Science and useful Arts, by securing for limited  
9 Times to Authors and Inventors the exclusive Right to their respective Writings and Discoveries.”  
10 U.S. Const. Art. I, Section 8. “Copyright law encourages people to create original works and thereby  
11 ‘ultimately serves the purpose of enriching the general public through access to creative works.’”  
12 *Fogerty v. Fantasy, Inc.*, 510 U.S. 517, 526 (1994).

13 106. The Supreme Court of the United States held that by “establishing a marketable right  
14 to the use of one’s expression, copyright supplies the economic incentive to create and disseminate  
15 ideas.” *Harper & Row Publisher, Inc. v. Nation Enters.*, 471 U.S. 539, 558 (1985).

16 107. The Copyright Act makes it illegal to publicly perform, display, distribute, or  
17 reproduce a copyrighted work except in limited instances, and provides for statutory damages,  
18 willful statutory damages, and the right to recover attorneys’ fees. 17 U.S.C. 501 *et seq.* The  
19 Copyright Act grants copyright owners the exclusive public display right, and control of the  
20 economic value of their protected works. This is true even where a copyrighted work is displayed  
21 somewhere online. Therefore, any person who downloads or even copies the work without consent  
22 is infringing on the owner’s exclusive rights to reproduction and/or distribution.

23 108. Defendant copied, downloaded, and otherwise misappropriated a vast trove of  
24 protected works available on the internet, including the exact digital version of Plaintiff Leovy’s  
25 copyrighted book, to develop the Gemini’s language model.

26 109. Defendant’s copying, storing processing, and reproducing of the entirety of Plaintiff  
27 Leovy’s copyrighted book and the copyrighted materials of the Class to train Gemini and other AI  
28 Products, infringed on Plaintiff Leovy’s and the Class Members’ exclusive rights in their



1 copyrighted works. Similarly, Defendant's blatant copying and unlawful appropriation of  
2 copyrighted works of others – images, books, song, etc. – infringed on Class Members' exclusive  
3 rights.

4 110. Also, after being trained (by illegally infringing on the copyrighted materials), the AI  
5 models are subject to fine tuning, wherein Google continues to train/re-train its AI models to better  
6 mimic the works, content, expressive style, protected designs, and other important parts of authors'  
7 works.

8 111. Plaintiff Leovy is the exclusive owner of the registered copyright in her work under  
9 17 U.S.C. § 106; in fact, Plaintiff Leovy registered the copyright for her book on February 20, 2015.

10 112. As exclusive rights holder, only Plaintiff Leovy or those Plaintiff Leovy has  
11 authorized may reproduce (i.e. copy, download), or distribute her property. Neither Plaintiff Leovy  
12 nor any Class Members authorized Defendant to use their works or make copies of their works.

13 113. Defendant generates billions of dollars on its AI technology, Gemini, which was  
14 trained on copyrighted works and materials without consent or compensation. Without this mass  
15 infringement, Gemini would not exist.

16 114. By training its Products on the protected works of millions of authors, Defendant  
17 engaged in unauthorized use, distribution, and reproduction of the copyrighted materials.

18 115. Defendant's conduct herein is willful because it is aware that stealing works from the  
19 entire internet will undoubtedly result in infringement, especially where Defendant is copying the  
20 databased and websites that are known to contain pirated books and works. Also, Defendant is aware  
21 that the stolen millions of works were registered with the U.S/ Copyright Office, because the  
22 copyright symbol appears in at least some datasets used by Defendant (C-4) more than 200 million  
23 times.

24 116. Defendant made copies and engaged in an unauthorized use of Plaintiff Leovy and  
25 Class Members' work for training and development of Gemini (as well as other AI Products).  
26 Defendant's infringement of copyrighted works was knowing, willful, and intentional, and thus  
27 subjects Defendant to liability for statutory damages under Section 504(c)(2) of the Copyright Act  
28 of up to \$150,000 per infringement. Furthermore, Defendant has sufficient resources to verify

1 whether or not the works on which Gemini and other AI Products were trained on are protected  
2 under copyright law.

3 117. Alternatively, even if Defendant was unaware and had no reason to believe that its  
4 actions constituted copyright infringement, Plaintiff Leovy and Class Members are entitled to  
5 \$200.00/per infringement.

6 118. As a direct and proximate cause of Defendant's conduct, Plaintiff Leovy and Class  
7 Members have suffered and will continue to suffer monetary damages in an amount to be determined  
8 at trial. Plaintiff Leovy and Class Members are entitled to statutory damages, actual damages,  
9 disgorgement of profits, injunctive and declaratory relief, and other remedies.

10 119. Because Plaintiff and the members of the proposed Class have been and continue to  
11 be irreparably injured due to Defendant's infringement and conduct described herein, for which no  
12 adequate remedy is available at law, Plaintiff and the Class are entitled to injunctive relief. Without  
13 permanent injunctive relief, Defendant will continue to infringe on the exclusive rights of Plaintiff  
14 and the proposed Class, unless its infringing activity is enjoined by this Court.

### 15 **PRAYER FOR RELIEF**

16 WHEREFORE, Plaintiff on behalf of herself and the Proposed Class that she seeks to  
17 represent, respectfully requests the following relief:

- 18 A. Certify this action as a class action pursuant to Rule 23 of the Federal Rules of Civil  
19 Procedure;
- 20 B. Appoint Plaintiff to represent the Class;
- 21 C. Appoint undersigned counsel to represent the Class;
- 22 D. Award compensatory damages to Plaintiff and the Class against Defendant for all  
23 damages sustained as a result of Defendant's wrongdoing, in an amount to be proven  
24 at trial, including interest;
- 25 E. Award statutory (including treble damages, where appropriate) damages to Plaintiff  
26 and the Class against Defendant;
- 27 F. Award nominal damages to Plaintiff and the Class against Defendant;
- 28 G. Non-restitutionary disgorgement of all profits that were derived, in whole or in part,

1 from Defendant's conduct;

2 H. Award punitive damages to Plaintiff and the Class against Defendant;

3 I. Permanently restrain Defendant, and its officers, agents, servants, employees, and  
4 attorneys, from the conduct at issue in this Action and otherwise violating its policies  
5 with consumers, and award all other appropriate injunctive and equitable relief  
6 deemed just and proper;

7 J. Award Plaintiff and the Class their reasonable costs and expenses incurred in this  
8 Action, including attorneys' fees, costs, and expenses; and

9 K. Grant Plaintiff and the Class such further relief as the Court deems appropriate.

10 **JURY TRIAL DEMANDED**

11 Plaintiff demands a jury trial on all triable issues.

12  
13 DATED: June 27, 2024

**CLARKSON LAW FIRM, P.C.**

14 */s/ Ryan J. Clarkson* \_\_\_\_\_

15 Ryan Clarkson, Esq.

16 Yana Hart, Esq.

17 Tracey Cowan, Esq.

18 Tiara Avanness, Esq.

19 *Counsel for Plaintiff and the Proposed Class*