Joseph R. Saveri (State Bar No. 130064)
**JOSEPH SAVERI LAW FIRM, LLP**
601 California Street, Suite 1000
San Francisco, CA 94108
Telephone:     (415) 500-6800
Facsimile:     (415) 395-9940
Email:          jsaveri@saverilawfirm.com

Matthew Butterick (State Bar No. 250953)
1920 Hillhurst Avenue, #406
Los Angeles, CA 90027
Telephone:     (323) 968-2632
Facsimile:     (415) 395-9940
Email:          mb@butoericklaw.com

Bryan L. Clobes (*pro hac vice*)
**CAFFERTY CLOBES MERIWETHER & SPRENGEL LLP**
205 N. Monroe Street
Media, PA 19063
Telephone:     (215) 864-2800
Email:          bclobes@caffertyclobes.com

*Counsel for Individual and Representative Plaintiffs*
*and the Proposed Class (continued on signature page)*

# UNITED STATES DISTRICT COURT
# NORTHERN DISTRICT OF CALIFORNIA
# SAN FRANCISCO DIVISION

| | |
|---|---|
| **IN RE OPENAI CHATGPT LITIGATION** | Master File No. 23-cv-3223-AMO |
| This document relates to:<br>23-cv-03223-AMO<br>23-cv-03416-AMO<br>23-cv-04625-AMO | **FIRST CONSOLIDATED AMENDED COMPLAINT**<br><br>**Class Action**<br><br>**Demand for Jury Trial** |

Plaintiffs Paul Tremblay, Sarah Silverman, Christopher Golden, Richard Kadrey, Ta-Nehisi Coates, Junot Díaz, Andrew Sean Greer, David Henry Hwang, Matthew Klam, Laura Lippman, Rachel Louise Snyder, and Jacqueline Woodson (collectively "Plaintiffs"), on behalf of themselves and all others similarly situated, bring this class-action complaint ("Complaint") against defendants OpenAI, Inc.; OpenAI, L.P.; OpenAI OpCo, L.L.C.; OpenAI GP, L.L.C.; OpenAI Startup Fund I, L.P.; OpenAI Startup Fund GP I, L.L.C.; and OpenAI Startup Fund Management, LLC (collectively, "OpenAI" or "Defendants"). Plaintiffs seek to recover injunctive relief and damages as a result of OpenAI's unlawful conduct.

## OVERVIEW

1. ChatGPT is a software product created, maintained, and sold by OpenAI.

2. ChatGPT is powered by two AI software programs called GPT-3.5 and GPT-4, also known as *large language models*. Rather than being programmed in the traditional way, a large language model is "trained" by copying massive amounts of text and extracting expressive information from it. This body of text is called the *training dataset*. Once a large language model has copied and ingested the text in its training dataset, it is able to emit convincingly naturalistic text outputs in response to user prompts.

3. A large language model's output is therefore entirely and uniquely reliant on the material in its training dataset. Every time it assembles a text output, the model relies on the information it extracted from its training dataset.

4. Plaintiffs and Class members are authors of books. Plaintiffs and Class members have registered copyrights in the books they published. Plaintiffs and Class members did not consent to the use of their copyrighted books as training material for ChatGPT. Nonetheless, their copyrighted materials were ingested and used to train ChatGPT.

5. Indeed, when ChatGPT is prompted, ChatGPT generates summaries of Plaintiffs' copyrighted works—something only possible if ChatGPT was trained on Plaintiffs' copyrighted works.

6. Defendants, by and through the use of ChatGPT, benefit commercially and profit significantly from the use of Plaintiffs' and Class members' copyrighted materials.

**JURISDICTION AND VENUE**

7.      This Court has subject-matter jurisdiction under 28 U.S.C. § 1331 because this case arises under the Copyright Act (17 U.S.C. § 101, *et seq.*).

8.      Jurisdiction and venue is proper in this judicial district under 28 U.S.C. § 1391(c)(2) because Defendant OpenAI, Inc. is headquartered in this District, and thus a substantial part of the events giving rise to the claims occurred in this District; and because a substantial part of the events giving rise to Plaintiffs' claims occurred in this District, and a substantial portion of the affected interstate trade and commerce was carried out in this District. Each Defendant has transacted business, maintained substantial contacts, and/or committed overt acts in furtherance of the illegal scheme and conspiracy throughout the United States, including in this District. Defendants' conduct has had the intended and foreseeable effect of causing injury to persons residing in, located in, or doing business throughout the United States, including in this District.

9.      Under Civil Local Rule 3-2(c), assignment of this case to the San Francisco Division is proper because this case pertains to intellectual-property rights, which is a district-wide case category under General Order No. 44, and therefore venue is proper in any courthouse in this District.

**PLAINTIFFS**

10.     Plaintiff Paul Tremblay is a writer who lives in Massachusetts and owns registered copyrights in multiple works, including *The Cabin at the End of the World*.

11.     Plaintiff Sarah Silverman is a writer and performer who lives in California and owns a registered copyright in one work, *The Bedwetter*.

12.     Plaintiff Christopher Golden is a writer who lives in Massachusetts and owns registered copyrights in multiple works, including *Ararat*.

13.     Plaintiff Richard Kadrey is a writer who lives in Pennsylvania and owns registered copyrights in multiple works, including *Sandman Slim*.

14.     Plaintiff Ta-Nehisi Coates is an author who lives in New York and owns registered copyrights in multiple works, including *The Water Dancer*.

FIRST CONSOLIDATED AMENDED COMPLAINT — 23-cv-3223-AMO

15.     Plaintiff Junot Díaz is an author who lives in Massachusetts and owns registered copyrights in multiple works, including *The Brief Wondrous Life of Oscar Wao*.

16.     Plaintiff Andrew Sean Greer is an author who lives in California and owns registered copyrights in multiple works, including *The Confessions of Max Tivoli*.

17.     Plaintiff David Henry Hwang is a playwright and screenwriter who lives in New York and owns registered copyrights in multiple works, including *The Dance and the Railroad*.

18.     Plaintiff Matthew Klam is an author who lives in Washington, D.C. and owns registered copyrights in multiple works, including *Who is Rich?*

19.     Plaintiff Laura Lippman is an author who lives in Maryland and owns registered copyrights in multiple works, including *What the Dead Know*.

20.     Plaintiff Rachel Louise Snyder is an author who lives in Washington, D.C. and owns registered copyrights in multiple works, including *What We've Lost is Nothing*.

21.     Plaintiff Jacqueline Woodson is an author who lives in New York and owns registered copyrights in multiple works, including *Brown Girl Dreaming*.

22.     A non-exhaustive list of copyright registrations owned by Plaintiffs is shown in Exhibit A. Together, these works are also referred to as the **Infringed Works**.

## DEFENDANTS

23.     Defendant OpenAI, Inc. is a Delaware nonprofit corporation with its principal place of business located at 3180 18th Street, San Francisco, CA 94110.

24.     Defendant OpenAI, L.P. is a Delaware limited partnership with its principal place of business located at 3180 18th Street, San Francisco, CA 94110. OpenAI, L.P. is a wholly owned subsidiary of OpenAI Inc. that is operated for profit. OpenAI, Inc. controls OpenAI, L.P. directly and through the other OpenAI entities.

25.     Defendant OpenAI OpCo, L.L.C. is a Delaware limited liability company with its principal place of business located at 3180 18th Street, San Francisco, CA 94110. OpenAI OpCo, L.L.C. is a wholly owned subsidiary of OpenAI, Inc. that is operated for profit. OpenAI, Inc. controls OpenAI OpCo, L.L.C. directly and through the other OpenAI entities.

FIRST CONSOLIDATED AMENDED COMPLAINT — 23-cv-3223-AMO

1    26.    Defendant OpenAI GP, L.L.C. ("OpenAI GP") is a Delaware limited liability company

2    with its principal place of business located at 3180 18th Street, San Francisco, CA 94110. OpenAI GP is

3    the general partner of OpenAI, L.P. OpenAI GP manages and operates the day-to-day business and

4    affairs of OpenAI, L.P. OpenAI GP was aware of the unlawful conduct alleged herein and exercised

5    control over OpenAI, L.P. throughout the Class Period. OpenAI, Inc. directly controls OpenAI GP.

6    27.    Defendant OpenAI Startup Fund I, L.P. ("OpenAI Startup Fund I") is a Delaware

7    limited partnership with its principal place of business located at 3180 18th Street, San Francisco, CA

8    94110. OpenAI Startup Fund I was instrumental in the foundation of OpenAI, L.P., including the

9    creation of its business strategy and providing initial funding. OpenAI Startup Fund I was aware of the

10   unlawful conduct alleged herein and exercised control over OpenAI, L.P. throughout the Class Period.

11   28.    Defendant OpenAI Startup Fund GP I, L.L.C. ("OpenAI Startup Fund GP I") is a

12   Delaware limited liability company with its principal place of business located at 3180 18th Street, San

13   Francisco, CA 94110. OpenAI Startup Fund GP I is the general partner of OpenAI Startup Fund I.

14   OpenAI Startup Fund GP I is a party to the unlawful conduct alleged herein. OpenAI Startup Fund GP

15   I manages and operates the day-to-day business and affairs of OpenAI Startup Fund I.

16   29.    Defendant OpenAI Startup Fund Management, LLC ("OpenAI Startup Fund

17   Management") is a Delaware limited liability company with its principal place of business located at

18   3180 18th Street, San Francisco, CA 94110. OpenAI Startup Fund Management is a party to the

19   unlawful conduct alleged herein. OpenAI Startup Fund Management was aware of the unlawful

20   conduct alleged herein and exercised control over OpenAI, L.P. throughout the Class Period.

21

22                              **AGENTS AND CO-CONSPIRATORS**

23   30.    The unlawful acts alleged against the Defendants in this class action complaint were

24   authorized, ordered, or performed by the Defendants' respective officers, agents, employees,

25   representatives, or shareholders while actively engaged in the management, direction, or control of the

26   Defendants' businesses or affairs. The Defendants' agents operated under the explicit and apparent

27   authority of their principals. Each Defendant, and its subsidiaries, affiliates, and agents operated as a

28   single unified entity.

FIRST CONSOLIDATED AMENDED COMPLAINT — 23-cv-3223-AMO

31.     Various persons and/or firms not named as Defendants may have participated as co-conspirators in the violations alleged herein and may have performed acts and made statements in furtherance thereof. Each acted as the principal, agent, or joint venturer of, or for, other Defendants with respect to the acts, violations, and common course of conduct alleged herein.

## FACTUAL ALLEGATIONS

32.     OpenAI creates and sells artificial-intelligence software products. *Artificial intelligence* is commonly abbreviated as "AI." AI software is designed to algorithmically simulate human reasoning or inference, often using statistical methods.

33.     Certain AI products created and sold by OpenAI are known as *large language models.* A large language model (or "LLM" for short) is AI software designed to parse and emit natural language. Though a large language model is a software program, it is not created the way most software programs are—that is, by human software engineers writing code. Rather, a large language model is "trained" by copying massive amounts of text from various sources and feeding these copies into the model. This corpus of input material is called the *training dataset*. During training, the large language model copies each piece of text in the training dataset and extracts expressive information from it. The large language model progressively adjusts its output to more closely resemble the sequences of words copied from the training dataset. Once the large language model has copied and ingested all this text, it is able to emit convincing simulations of natural written language as it appears in the training dataset.

34.     Much of the material in OpenAI's training datasets, however, comes from copyrighted works—including books written by Plaintiffs—that were copied by OpenAI without consent, without credit, and without compensation.

35.     Authors, including Plaintiffs, publish books with certain copyright management information. This information includes the book's title, the ISBN number or copyright number, the author's name, the copyright holder's name, and terms and conditions of use. This information is most commonly found on the back of the book's title page and is customarily included in all books, regardless of genre.

36.     OpenAI made a series of large language models, including without limitation GPT-1 (released June 2018), GPT-2 (February 2019), GPT-3 (May 2020), GPT-3.5 (March 2022), and most recently GPT-4 (March 2023). "GPT" is an abbreviation for "generative pre-trained transformer," where *pre-trained* refers to the use of textual material for training, *generative* refers to the model's ability to emit text, and *transformer* refers to the underlying training algorithm. OpenAI offers certain language models in variant forms: for instance, the GPT-4 family of models includes publicly accessible variants called 'gpt-4-0125-preview,' 'gpt-4-turbo-preview,' and 'gpt-4-32k;' the GPT-3.5 Turbo family of models includes publicly accessible variants called 'gpt-3.5-turbo-0125,' 'gpt-3.5-turbo-1106,' and 'gpt-3.5-turbo-instruct.' On information and belief, OpenAI has made other language-model variants that are in commercial use but are not publicly accessible. In an interview with the Financial Times in November 2023, OpenAI CEO Sam Altman confirmed that GPT-5 is under development. Together, OpenAI's large language models, including any in development, will be referred to as the "OpenAI Language Models."[1]

37.     Many kinds of material have been used to train large language models. Books, however, have always been a key ingredient in training datasets for large language models because books offer the best examples of high-quality longform writing.

38.     For instance, in its June 2018 paper introducing GPT-1 (called "Improving Language Understanding by Generative Pre-Training"), OpenAI revealed that it trained GPT-1 on BookCorpus, a collection of "over 7,000 unique unpublished books from a variety of genres including Adventure, Fantasy, and Romance." OpenAI confirmed why a dataset of books was so valuable: "Crucially, it contains long stretches of contiguous text, which allows the generative model to learn to condition on long-range information." Hundreds of large language models have been trained on BookCorpus, including those made by OpenAI, Google, Amazon, and others.

39.     BookCorpus, however, is a controversial dataset. It was assembled in 2015 by a team of AI researchers for the purpose of training language models. They copied the books from a website called Smashwords.com that hosts unpublished novels that are available to readers at no cost. Those

---

[1] The definition of "OpenAI Language Models" encompasses any language models developed (or in development) by OpenAI, irrespective of whether those models underly ChatGPT.

FIRST CONSOLIDATED AMENDED COMPLAINT — 23-cv-3223-AMO

1 | novels, however, are largely under copyright. They were copied into the BookCorpus dataset without

2 | consent, credit, or compensation to the authors.

3 |     40.    OpenAI also copied many books while training GPT-3. In the July 2020 paper

4 | introducing GPT-3 (called "Language Models are Few-Shot Learners"), OpenAI disclosed that 15% of

5 | the enormous GPT-3 training dataset came from "two internet-based books corpora" that OpenAI

6 | simply called "Books1" and "Books2."

7 |     41.    Tellingly, OpenAI has never revealed what books are part of the Books1 and Books2

8 | datasets—though there are some clues. First, OpenAI admitted these are "internet-based books

9 | corpora." Second, both Books1 and Books2 are apparently much larger than BookCorpus. Based on

10 | numbers included in OpenAI's paper about GPT-3, Books1 is apparently about nine times larger;

11 | Books2 is about 42 times larger. Since BookCorpus contained about 7,000 titles, this suggests Books1

12 | would contain about 63,000 titles; Books2 would contain about 294,000 titles.

13 |     42.    But there are only a handful of "internet-based books corpora" that would be able to

14 | deliver this much material.

15 |     43.    As noted in ¶ 41, the OpenAI Books1 dataset can be estimated to contain about 63,000

16 | titles. Project Gutenberg is an online archive of e-books whose copyright has expired. In September

17 | 2020, Project Gutenberg claimed to have "over 60,000" titles. Project Gutenberg has long been

18 | popular for training AI systems due to the lack of copyright. In 2018, a team of AI researchers created

19 | the "Standardized Project Gutenberg Corpus," which contained "more than 50,000 books." On

20 | information and belief, the OpenAI Books1 dataset is based on either the Standardized Project

21 | Gutenberg Corpus or Project Gutenberg itself, because of the roughly similar sizes of the two datasets.

22 |     44.    As noted in ¶ 41, the OpenAI Books2 dataset can be estimated to contain about 294,000

23 | titles. The only "internet-based books corpora" that have ever offered that much material are notorious

24 | "shadow library" websites like Library Genesis (aka LibGen), Z-Library (aka B-ok), Sci-Hub, and

25 | Bibliotik. The books aggregated by these websites have also been available in bulk via torrent systems.

26 | These flagrantly illegal shadow libraries have long been of interest to the AI-training community: for

27 | instance, an AI training dataset published in December 2020 by EleutherAI called "Books3" includes a

28 | recreation of the Bibliotik collection and contains nearly 200,000 books. On information and belief, the

7

1  OpenAI Books2 dataset includes books copied from these "shadow libraries," because those are the

2  largest sources of trainable books most similar in nature and size to OpenAI's description of Books2.

3        45.     In March 2023, OpenAI's paper introducing GPT-4 contained no information about its

4  dataset at all: OpenAI claimed that "[g]iven both the competitive landscape and the safety implications

5  of large-scale models like GPT-4, this report contains no further details about … dataset construction."

6  Later in the paper, OpenAI concedes it did "filter[ ] our dataset … to specifically reduce the quantity of

7  inappropriate erotic text content."

8

9  **INTERROGATING THE OPENAI LANGUAGE MODELS USING CHATGPT**

10        46.     ChatGPT is a language model created and sold by OpenAI. As its name suggests,

11  ChatGPT is designed to offer a conversational style of interaction with a user. OpenAI offers ChatGPT

12  through a web interface to individual users for $20 per month. Through the web interface, users can

13  choose to use two versions of ChatGPT: one based on the GPT-3.5 model, and one based on the newer

14  GPT-4 model.

15        47.     OpenAI also offers ChatGPT to software developers through an application-

16  programming interface (or "API"). The API allows developers to write programs that exchange data

17  with ChatGPT. Access to ChatGPT through the API is billed on the basis of usage.

18        48.     Regardless of how it is accessed—either through the web interface or through the API—

19  ChatGPT allows users to enter text prompts, which ChatGPT then attempts to respond to in a natural

20  way, i.e., ChatGPT can generate answers in a coherent and fluent way that closely mimics human

21  language. If a user prompts ChatGPT with a question, ChatGPT will answer. If a user prompts

22  ChatGPT with a command, ChatGPT will obey. If a user prompts ChatGPT to summarize a

23  copyrighted book, it will do so.

24        49.     ChatGPT's output, like other LLMs, relies on the data upon which it is trained to

25  generate new content. LLMs generate output based on patterns and connections drawn from the

26  training data. For example, if an LLM is prompted to generate a writing in the style of a certain author,

27  the LLM would generate content based on patterns and connections it learned from analysis of that

28  author's work within its training data.

50.   On information and belief, the reason ChatGPT can accurately summarize a certain copyrighted book is because that book was copied by OpenAI and ingested by the underlying OpenAI Language Model (either GPT-3.5 or GPT-4) as part of its training data.

51.   When ChatGPT was prompted to summarize books written by each of the Plaintiffs, it generated very accurate summaries. These summaries are attached as **Exhibit B**. The summaries get some details wrong, which is expected, since a large language model mixes together expressive material derived from many sources. Still, the rest of the summaries are accurate, which means that ChatGPT retains knowledge of particular works in the training dataset and is able to output similar textual content.

## CLASS ALLEGATIONS

52.   The "**Class Period**" as defined in this Complaint begins on at least June 28, 2020 and runs through the present. Because Plaintiffs do not yet know when the unlawful conduct alleged herein began, but believe, on information and belief, that the conduct likely began earlier than June 28, 2020, Plaintiffs reserve the right to amend the Class Period to comport with the facts and evidence uncovered during further investigation or through discovery.

53.   **Class definition**. Plaintiffs bring this action for damages and injunctive relief as a class action under Federal Rules of Civil Procedure 23(a), 23(b)(2), and 23(b)(3), on behalf of the following Class:

> **All persons or entities domiciled in the United States that own a United States copyright in any work that was used as training data for the OpenAI Language Models during the Class Period.**

54.   This Class definition excludes:

a.   any of the Defendants named herein;

b.   any of the Defendants' co-conspirators;

c.   any of Defendants' parent companies, subsidiaries, and affiliates;

d.   any of Defendants' officers, directors, management, employees, subsidiaries, affiliates, or agents;

FIRST CONSOLIDATED AMENDED COMPLAINT — 23-cv-3223-AMO

1            e.     all governmental entities; and

2            f.     the judges and chambers staff in this case, as well as any members of their

3                immediate families.

4       55.     **Numerosity**. Plaintiffs do not know the exact number of members in the Class. This

5 information is in the exclusive control of Defendants. On information and belief, there are at least

6 thousands of members in the Class geographically dispersed throughout the United States. Therefore,

7 joinder of all members of the Class in the prosecution of this action is impracticable.

8       56.     **Typicality**. Plaintiffs' claims are typical of the claims of other members of the Class

9 because Plaintiffs and all members of the Class were damaged by the same wrongful conduct of

10 Defendants as alleged herein, and the relief sought herein is common to all members of the Class.

11       57.     **Adequacy**. Plaintiffs will fairly and adequately represent the interests of the members of

12 the Class because the Plaintiffs have experienced the same harms as the members of the Class and have

13 no conflicts with any other members of the Class. Furthermore, Plaintiffs retained and are represented

14 by sophisticated and competent counsel who are experienced in prosecuting federal and state class

15 actions, as well as other complex litigation.

16       58.     **Commonality and predominance**. Numerous questions of law or fact common to each

17 Class arise from Defendants' conduct:

18           a.   whether Defendants violated the copyrights of Plaintiffs and the Class when they

19              downloaded copies of Plaintiffs' copyrighted books and used them to train ChatGPT;

20           b.   whether ChatGPT itself is an infringing derivative work based on Plaintiffs' copyrighted

21              books;

22           c.   Whether Defendants' conduct alleged herein constitutes Unfair Competition under

23              California Business and Professions Code § 17200 *et seq*.

24           d.   Whether this Court should enjoin Defendants from engaging in the unlawful conduct

25              alleged herein. And what the scope of that injunction would be.

26           e.   Whether any affirmative defense excuses Defendants' conduct.

27           f.   Whether any statutes of limitation constrain the potential recovery for Plaintiffs and the

28              Class.

FIRST CONSOLIDATED AMENDED COMPLAINT — 23-cv-3223-AMO

59.     These and other questions of law and fact are common to the Class and predominate over any questions affecting the members of the Class individually.

60.     **Other class considerations**. Defendants acted on grounds generally applicable to the Class. This class action is superior to alternatives, if any, for the fair and efficient adjudication of this controversy. Prosecuting the claims pleaded herein as a class action will eliminate the possibility of repetitive litigation and inconsistent results. There will be no material difficulty in the management of this Action as a class action. Further, final injunctive relief is appropriate with respect to the Class as a whole.

61.     The prosecution of separate actions by individual Class members would create the risk of inconsistent or varying adjudications, establishing incompatible standards of conduct for Defendants.

## COUNT 1

### DIRECT COPYRIGHT INFRINGEMENT

### 17 U.S.C. § 501

62.     Plaintiffs incorporate by reference the preceding factual allegations.

63.     As the owners of the registered copyrights in books used to train the OpenAI Language Models, Plaintiffs hold the exclusive rights to those texts under 17 U.S.C. § 106.

64.     Plaintiffs never authorized OpenAI to make copies of their books, make derivative works, publicly display copies (or derivative works), or distribute copies (or derivative works). All those rights belong exclusively to Plaintiffs under copyright law.

65.     On information and belief, to train the OpenAI Language Models, OpenAI relied on harvesting mass quantities of textual material from the public internet, including Plaintiffs' books, which are available in digital formats.

66.     OpenAI made copies of Plaintiffs' books during the training process of the OpenAI Language Models without Plaintiffs' permission. Specifically, OpenAI copied at least the Infringed Works in Exhibit A.

67.     Licensing copyrighted material to train AI models is plainly feasible. It already happens. Indeed, OpenAI itself sought to license copyrighted materials to train its LLMs. For instance, OpenAI reached agreements with the Associated Press and Axel Springer to license textual materials for the purpose of LLM training. OpenAI has reportedly been in negotiations with other publishers as well.

68.     Because the OpenAI Language Models cannot function without the expressive information extracted from Plaintiffs' works (and others) and retained inside them, the OpenAI Language Models are themselves infringing derivative works, made without Plaintiffs' permission and in violation of their exclusive rights under the Copyright Act.

69.     Plaintiffs have been injured by OpenAI's acts of direct copyright infringement. Plaintiffs are entitled to statutory damages, actual damages, restitution of profits, and other remedies provided by law.

## COUNT 2

### UCL — Unfair Competition

### Cal. Bus. & Prof. Code §§ 17200 et seq.

70.     Plaintiffs incorporate by reference the preceding factual allegations.

71.     Defendants engaged in unfair business practices by, among other things, using Plaintiffs' Infringed Works to train ChatGPT without permission from Plaintiffs or Class members.

72.     The unfair business practices described herein violate California Business and Professions Code § 17200 *et seq.* (the "UCL").

73.     The unfair business practices described herein violate the UCL because they are unfair, immoral, unethical, oppressive, unscrupulous, or injurious to consumers, and because Defendants used Plaintiffs' protected works to train ChatGPT for Defendants' own commercial profit without the authorization of Plaintiffs or the Class. Defendants unfairly profit from and take credit for developing a commercial product based on unattributed reproductions of those stolen writings and ideas.

74.     The unlawful business practices described herein violate the UCL because consumers are likely to be deceived. Defendants knowingly and secretively trained ChatGPT using unauthorized

copies of Plaintiffs' copyrighted work. Defendants deceptively marketed their product in a manner that fails to attribute the success of their product to the copyrighted work on which it is based.

## DEMAND FOR JUDGMENT

WHEREFORE, Plaintiffs request that the Court enter judgment on their behalf and on behalf of the Class defined herein, by ordering:

a) This Action may proceed as a class action, with Plaintiffs serving as Class Representatives, and with Plaintiffs' counsel as Class Counsel.

b) A declaration that Defendants have infringed Plaintiffs' and the Class' exclusive copyrights in the Infringed Works under the Copyright Act.

c) A declaration that such infringement is willful.

d) Judgment in favor of Plaintiffs and the Class and against Defendants.

e) An award of statutory and other damages under 17 U.S.C. § 504 for Defendants' willful infringement of Plaintiffs' and the Class' exclusive copyrights in the Infringed Works.

f) Reasonable attorneys' fees and costs as available under 17 U.S.C. § 505 or other applicable statute.

g) Pre- and post-judgment interest on the damages awarded to Plaintiffs and the Class, and that such interest be awarded at the highest legal rate from and after the date this class action complaint is first served on Defendants.

h) Defendants are to be jointly and severally responsible financially for the costs and expenses of a Court-approved notice program through post and media designed to give immediate notification to the Class.

i) Further relief for Plaintiffs and the Class as may be just and proper.

1

**JURY TRIAL DEMANDED**

2      Under Federal Rule of Civil Procedure 38(b), Plaintiffs demand a trial by jury of all the claims

3  asserted in this Complaint so triable.

4

5  Dated: March 13, 2024                    By:      _/s/ *Joseph R. Saveri*_

6                                                     Joseph R. Saveri

7  Joseph R. Saveri (State Bar No. 130064)       Bryan L. Clobes (*pro hac vice*)
   Cadio Zirpoli (State Bar No. 179108)          **CAFFERTY CLOBES MERIWETHER**
8  Christopher K.L. Young (State Bar No. 318371)  **& SPRENGEL LLP**
   Holden Benon (State Bar No. 325847)           205 N. Monroe Street
9  Aaron Cera (State Bar No. 351163)             Media, PA 19063
10 **JOSEPH SAVERI LAW FIRM, LLP**             Telephone:    (215) 864-2800
   601 California Street, Suite 1000             Email:        bclobes@caffertyclobes.com
11 San Francisco, California 94108
   Telephone:    (415) 500-6800                  Alexander J. Sweatman (*pro hac vice* forthcoming)
12 Facsimile:    (415) 395-9940                  **CAFFERTY CLOBES MERIWETHER**
13 Email:        jsaveri@saverilawfirm.com       **& SPRENGEL LLP**
                 czirpoli@saverilawfirm.com      135 South LaSalle Street, Suite 3210
14               cyoung@saverilawfirm.com        Chicago, IL 60603
                 hbenon@saverilawfirm.com        Telephone:    (312) 782-4880
15               acera@saverilawfirm.com         Email:        asweatman@caffertyclobes.com
16

17 Matthew Butterick (State Bar No. 250953)       *Counsel for Individual and Representative*
   1920 Hillhurst Avenue, #406                    *Plaintiffs and the Proposed Class*
18 Los Angeles, CA 90027
19 Telephone:    (323) 968-2632
   Facsimile:    (415) 395-9940
20 Email:        mb@butticklaw.com

21

22

23

24

25

26

27

28