

BURSOR & FISHER, P.A.

L. Timothy Fisher (State Bar No. 191626)
Joshua B. Glatt (State Bar No. 354064)
1990 North California Blvd., 9th Floor
Walnut Creek, CA 94596
Telephone: (925) 300-4455
Facsimile: (925) 407-2700
E-mail: ltfisher@bursor.com
jglatt@bursor.com

BURSOR & FISHER, P.A.

Joseph I. Marchese (*pro hac vice* forthcoming)
Julian C. Diamond (*pro hac vice* forthcoming)
1330 Avenue of the Americas, 32nd Floor
New York, NY 10019
Telephone: (646) 837-7150
Facsimile: (212) 989-9163
E-Mail: jmarchese@bursor.com
jdiamond@bursor.com

Attorneys for Plaintiff

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

DAVID MILLETTE, individually and on behalf
of all others similarly situated,

Plaintiff,

v.

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

Defendants.

Case No.

CLASS ACTION COMPLAINT

JURY TRIAL DEMANDED

1 Plaintiff David Millette, (hereinafter “Plaintiff”), brings this action on behalf of himself and
2 all others similarly situated against Defendants OpenAI, Inc.; OpenAI, L.P.; OpenAI OpCo,
3 L.L.C.; OpenAI GP, L.L.C.; OpenAI Startup Fund I, L.P.; OpenAI Startup Fund GP I, L.L.C.; and
4 OpenAI Startup Fund Management, LLC (collectively, “OpenAI” or “Defendants”). Plaintiff
5 seeks to recover injunctive relief and damages as a result of Defendants’ unlawful conduct.
6 Plaintiff makes the following allegations pursuant to the investigation of his counsel and are based
7 upon information and belief, except as to the allegations specifically pertaining to himself, which
8 are based on personal knowledge.

9 NATURE OF THE CASE

- 10 1. ChatGPT is a software product created, maintained, and sold by OpenAI.
- 11 2. ChatGPT is currently powered by artificial-intelligence (hereinafter “AI”) software
12 programs called GPT-3.5, GPT-4, and GPT-4o also known as *large language models*. A large
13 language model is “trained” by copying massive amounts of text and extracting expressive
14 information from it. This body of text is called the *training dataset*. Once a large language model
15 has copied and ingested the text in its training dataset, it is able to emit convincingly naturalistic
16 text outputs in response to user prompts.
- 17 3. Large language models’ output is therefore entirely and uniquely reliant on the
18 material in their training dataset. Every time they assemble text, video or image outputs, the
19 models rely on the information they extracted from their training dataset.
- 20 4. This case addresses the surreptitious, non-consensual transcription of millions of
21 YouTube users’ videos by Defendants to train Defendants’ AI software products. For years,
22 YouTube has been a popular video sharing platform that allows content creators and users to
23 upload and share videos with audiences worldwide. However, unbeknownst to those who upload
24 videos to YouTube, Defendants have been covertly transcribing YouTube videos to create training
25 datasets that they then use to train their AI products.
- 26 5. Plaintiff and Class members are YouTube users and video creators. Plaintiff and
27 Class members have retained ownership rights in their uploaded videos, per YouTube’s Terms of
28

1 Service. Plaintiff and Class members did not consent to the use of their videos as training material
2 for ChatGPT. Nonetheless, their materials were transcribed and used to train ChatGPT.

3 6. By transcribing and using these videos in this way, Defendants profit from
4 Plaintiff's and class members' data time and time again. As Defendants' AI products become more
5 sophisticated through the use of training datasets, they become more valuable to prospective and
6 current users, who purchase subscriptions to access Defendants' AI products.

7 7. By collecting and using this data without consent, Defendants have profited
8 significantly from the use of Plaintiff's and Class members' materials, violated California's Unfair
9 Competition Law ("UCL"), and been unjustly enriched at Plaintiff and Class members' expense.

10 JURISDICTION AND VENUE

11 8. This Court has subject matter Jurisdiction over this action pursuant to the Class
12 Action Fairness Act ("CAFA"), 28 U.S.C. § 1332(d)(2) because this is a class action in which at
13 least one member of the class is a citizen of a state different from any Defendants, the amount in
14 controversy exceeds \$5 million, exclusive of interest and costs, and the proposed class contains
15 more than 100 members.

16 9. This Court has personal jurisdiction over the Defendants because Defendants
17 maintain their principal places of business in this District and because a substantial part of the
18 events or omissions giving rise to the claims asserted herein occurred in this District.

19 10. Venue is proper in this district pursuant to 28 U.S.C. § 1391 because a substantial
20 part of the events or omissions giving rise to the claims asserted herein occurred in this District and
21 because Defendants maintain their principal places of business in this District.

22 PARTIES

23 11. Plaintiff David Millette is a resident of Douglas, Massachusetts. Plaintiff created a
24 YouTube account in or around June 2009. During that entire time, Plaintiff has retained ownership
25 rights to the video content he has uploaded to YouTube, per YouTube's Terms of Service.

26 12. Defendants transcribed Plaintiff's videos to train their AI software products.

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

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

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

9 16. Defendant OpenAI GP, L.L.C. (“OpenAI GP”) is a Delaware limited liability
10 company with its principal place of business located at 3180 18th Street, San Francisco, CA 94110.
11 OpenAI GP is the general partner of OpenAI, L.P. OpenAI GP manages and operates the day-to-
12 day business and affairs of OpenAI, L.P. OpenAI GP was aware of the unlawful conduct alleged
13 herein and exercised control over OpenAI, L.P. throughout the Class Period. OpenAI, Inc. directly
14 controls OpenAI GP.

15 17. Defendant OpenAI Startup Fund I, L.P. (“OpenAI Startup Fund I”) is a Delaware
16 limited partnership with its principal place of business located at 3180 18th Street, San Francisco,
17 CA 94110. OpenAI Startup Fund I was instrumental in the foundation of OpenAI, L.P., including
18 the creation of its business strategy and providing initial funding. OpenAI Startup Fund I was
19 aware of the unlawful conduct alleged herein and exercised control over OpenAI, L.P. throughout
20 the Class Period.

21 18. Defendant OpenAI Startup Fund GP I, L.L.C. (“OpenAI Startup Fund GP I”) is a
22 Delaware limited liability company with its principal place of business located at 3180 18th Street,
23 San Francisco, CA 94110. OpenAI Startup Fund GP I is the general partner of OpenAI Startup
24 Fund I. OpenAI Startup Fund GP I is a party to the unlawful conduct alleged herein. OpenAI
25 Startup Fund GP I manages and operates the day-to-day business and affairs of OpenAI Startup
26 Fund I.

27 19. Defendant OpenAI Startup Fund Management, LLC (“OpenAI Startup Fund
28 Management”) is a Delaware limited liability company with its principal place of business located

1 at 3180 18th Street, San Francisco, CA 94110. OpenAI Startup Fund Management is a party to the
2 unlawful conduct alleged herein. OpenAI Startup Fund Management was aware of the unlawful
3 conduct alleged herein and exercised control over OpenAI, L.P. throughout the Class Period.

4 20. Each of the Defendants acted jointly to perpetrate the acts described herein. At all
5 times relevant to the allegations in this matter, each of these Defendants acted in concert with, with
6 the knowledge and approval of, and/or as the agent of the other Defendants within the course and
7 scope of the agency, regarding the acts and omissions alleged.

8 **GENERAL BACKGROUND**

9 21. OpenAI creates and sells artificial-intelligence (AI) software products. AI software
10 is designed to algorithmically simulate human reasoning or inference, often using statistical
11 methods.

12 22. Certain AI products created and sold by OpenAI are known as large language
13 models. Large language models (“LLMs”) are types of AI software designed to parse and emit
14 natural language. Though LLMs are software programs, they are not created the way most
15 software programs are—that is, by human software engineers writing code. Rather, LLMs are
16 “trained” by copying massive amounts of text from various sources and feeding these copies into
17 the model. During training, these models copy each piece of information in the training dataset and
18 extract expressive information from it. The LLMs progressively adjust their output to more closely
19 resemble the images, videos, and sequences of words copied from the training dataset. Once these
20 models have copied and ingested all these inputs, they are able to emit convincing simulations of
21 natural written language, as well as videos and images as they appear in the training dataset.

22 23. Much of the material in OpenAI’s training datasets, however, comes from works—
23 including videos created and uploaded by Plaintiff—that were copied by OpenAI without consent,
24 without credit, and without compensation.

25 24. OpenAI made a series of large language models, including without limitation GPT-1
26 (released June 2018), GPT-2 (February 2019), GPT-3 (May 2020), GPT-3.5 (March 2022), GPT-4
27 (March 2023) and GPT-4o (May 2024). “GPT” is an abbreviation for “generative pre-trained
28 transformer,” where pre-trained refers to the use of textual material for training, generative refers to

1 the model’s ability to emit text, and transformer refers to the underlying training algorithm.
2 OpenAI offers certain language models in variant forms: for instance, the GPT-4 family of models
3 includes publicly accessible variants called ‘gpt-4-0125-preview,’ ‘gpt-4-turbo-preview,’ and ‘gpt-
4 4-32k;’ the GPT-3.5 Turbo family of models includes publicly accessible variants called ‘gpt-3.5-
5 turbo-0125,’ ‘gpt-3.5-turbo-1106,’ and ‘gpt-3.5-turbo-instruct.’ On information and belief,
6 OpenAI has made other language-model variants that are in commercial use but are not publicly
7 accessible. In an interview with the Financial Times in November 2023, OpenAI CEO Sam
8 Altman confirmed that GPT-5 is under development. Together, OpenAI’s large language models,
9 including any in development, will be referred to as the “OpenAI Language Models.”¹

10 25. Many kinds of material have been used to train large language models. Video
11 transcriptions, however, are a key ingredient in training datasets for large language models because
12 they offer copious examples of natural language.

13 26. In 2022, OpenAI released an automatic speech recognition (ASR) system called
14 Whisper. The Whisper model, which transcribes audio into text, was trained on 680,000 hours of
15 data collected from across the web. Tellingly, the exact names of the speech recognition corpora
16 on which Whisper was trained are unavailable. But one of the world’s largest open multilingual
17 speech corpora, VoxPopuli, contains only 400,000 hours of unlabeled speech data. Libriheavy, an
18 ASR corpus considered one of the largest freely available corpora of speech with supervisions,
19 only consists of 50,000 hours of English speech derived from LibriVox.

20 27. There are only a handful of publicly available, internet-based speech corpora that
21 can be utilized as training data for LLMs. As demonstrated, the two biggest corpora combined
22 (VoxPopuli and Libriheavy) still fall more than 200,000 hours short of the duration of speech that
23 comprises Whisper’s training dataset.

24 28. A New York Times report claims that Whisper is capable of transcribing the audio
25 from YouTube videos, and that an OpenAI team that included OpenAI’s president, Greg
26 Brockman, transcribed more than one million hours of video from YouTube.

27 _____
28 ¹ The definition of “OpenAI Language Models” encompasses any language models developed (or
in development) by OpenAI, irrespective of whether those models underly ChatGPT.

- 1 a. Whether Defendants violated the rights of Plaintiff and the Class when they
- 2 transcribed Plaintiff's videos and used those transcriptions as part of their AI
- 3 software's training datasets;
- 4 b. Whether Defendant OpenAI's conduct alleged herein constitutes Unfair
- 5 Competition under California Business and Professions Code § 17200 *et seq.*
- 6 c. Whether this Court should enjoin Defendants from engaging in the unlawful
- 7 conduct alleged herein, and what the scope of that injunction would be.
- 8 d. Whether any affirmative defense excuses Defendants' conduct.
- 9 e. Whether any statutes of limitation constrain the potential recovery for
- 10 Plaintiff and the Class.
- 11 f. Whether Plaintiff and the other Class members are entitled to restitution or
- 12 other relief.

13 36. **Typicality.** Plaintiff's claims are typical of the claims of the other members of the
14 Class in that, among other things, all Class members were similarly situated and were comparably
15 injured through Defendants' wrongful conduct as set forth herein. Further, there are no defenses
16 available to Defendants that are unique to Plaintiff.

17 37. **Adequacy of Representation.** Plaintiff will fairly and adequately protect the
18 interests of the Class. Plaintiff has retained counsel that is highly experienced in complex
19 consumer class action litigation, and Plaintiff intends to vigorously prosecute this action on behalf
20 of the Class. Furthermore, Plaintiff has no interests that are antagonistic to those of the Class.

21 38. **Superiority.** A class action is superior to all other available means for the fair and
22 efficient adjudication of this controversy. The damages or other financial detriment suffered by
23 individual Class members are relatively small compared to the burden and expense of individual
24 litigation of their claims against Defendants. It would thus be virtually impossible for the Class to
25 obtain effective redress for the wrongs committed against the members on an individual basis.
26 Furthermore, even if Class members could afford such individualized litigation, the court system
27 could not. Individualized litigation would create the danger of inconsistent or contradictory
28 judgments arising from the same set of facts. Individualized litigation would also increase the

1 delay and expense to all parties and the court system from the issues raised by this action. By
2 contrast, the class action device provides the benefits of adjudication of these issues in a single
3 proceeding, economies of scale, and comprehensive supervision by a single court, and presents no
4 unusual management difficulties under the circumstances.

5 39. Further, Defendants have acted and refused to act on grounds generally applicable
6 to the proposed Class, thereby making appropriate final injunctive and declaratory relief with
7 respect to the Class as a whole.

8 **CAUSES OF ACTION**

9 **COUNT I**

10 **Unjust Enrichment or Restitution
(On behalf of Plaintiff and the Class)**

11 40. Plaintiff incorporates by reference and re-alleges each and every allegation set forth
12 above as though fully set forth herein.

13 41. Plaintiff brings this claim individually and on behalf of members of the Class
14 against the Defendants.

15 42. To the extent required by law, Plaintiff brings this claim in the alternative to any
16 legal claims that may be alleged.

17 43. Plaintiff also alternatively alleges this claim as a Quasi-Contract or Non-Quasi-
18 Contract Claim for Restitution and Disgorgement.

19 44. Plaintiff and Class members unwittingly conferred a benefit upon Defendants.
20 OpenAI acquired valuable information from Plaintiff and Class members' videos to expand their
21 AI software's training datasets and used that information to develop and improve their products. In
22 using Plaintiff's information to refine their Language Models, OpenAI made their products more
23 valuable to prospective and current users, who purchase subscriptions to access them. Plaintiff and
24 Class members received nothing from this transaction. Plaintiff lacks an adequate remedy at law,
25 and pleads this cause of action in the alternative to the extent Plaintiff is required to do so.

26 45. Defendants have knowledge of such benefits.
27
28

- 1 g. For injunctive relief as the Court may deem proper; and
2 h. For an order awarding Plaintiff and the Class their reasonable attorneys' fees and
3 expenses and costs of suit.

4 **DEMAND FOR TRIAL BY JURY**

5 Pursuant to Federal Rule of Civil Procedure 38(b), Plaintiff demands a trial by jury of any
6 and all issues in this action so triable of right.

7
8 Dated: August 2, 2024

BURSOR & FISHER, P.A.

9 By: /s/ L. Timothy Fisher
10 L. Timothy Fisher

11 L Timothy Fisher (State Bar No. 191626)
12 Joshua B. Glatt (State Bar No. 354064)
13 1990 North California Blvd., 9th Floor
14 Walnut Creek, CA 94596
15 Telephone: (925) 300-4455
16 Facsimile: (925) 407-2700
17 E-mail: ltfisher@bursor.com
18 jglatt@bursor.com

BURSOR & FISHER, P.A.

19 Joseph I. Marchese (*pro hac vice* forthcoming)
20 Julian C. Diamond (*pro hac vice* forthcoming)
21 1330 Avenue of the Americas, 32nd Floor
22 New York, NY 10019
23 Telephone: (646) 837-7150
24 Facsimile: (212) 989-9163
25 E-Mail: jmarchese@bursor.com
26 jdiamond@bursor.com

27 *Attorneys for Plaintiff*
28