

1 Joseph R. Saveri (State Bar No. 130064)  
 2 Cadio Zirpoli (State Bar No. 179108)  
 3 Christopher K.L. Young (State Bar No. 318371)  
 4 Louis A. Kessler (State Bar No. 243703)  
 5 Elissa A. Buchanan (State Bar No. 249996)  
 6 Travis Manfredi (State Bar No. 281779)  
 7 William W. Castillo Guardado (State Bar No. 294159)  
 8 Holden J. Benon (State Bar No. 325847)  
**JOSEPH SAVERI LAW FIRM, LLP**  
 601 California Street, Suite 1000  
 San Francisco, California 94108  
 Telephone: (415) 500-6800  
 Facsimile: (415) 395-9940  
 Email: jsaveri@saverilawfirm.com  
 czirpoli@saverilawfirm.com  
 cyoung@saverilawfirm.com  
 lkessler@saverilawfirm.com  
 eabuchanan@saverilawfirm.com  
 tmanfredi@saverilawfirm.com  
 wcastillo@saverilawfirm.com  
 hbenon@saverilawfirm.com

*Counsel for Individual and Representative  
 Plaintiffs and the Proposed Class*

[Additional Counsel Listed on Signature Page]

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

J. DOE 1, J. DOE 2, J. DOE 3, J. DOE 4, and J. DOE 5,  
 individually and on behalf of all others similarly  
 situated,

*Individual and Representative Plaintiffs,*

v.

GITHUB, INC., a Delaware corporation;  
 MICROSOFT CORPORATION, a Washington  
 corporation; OPENAI, INC., a Delaware nonprofit  
 corporation; OPENAI, L.P., a Delaware limited  
 partnership; OPENAI OPCO, L.L.C., a Delaware  
 limited liability company; OPENAI GP, L.L.C., a  
 Delaware limited liability company; OPENAI  
 STARTUP FUND GP I, L.L.C., a Delaware limited  
 liability company; OPENAI STARTUP FUND I, L.P.,  
 a Delaware limited partnership; OPENAI STARTUP  
 FUND MANAGEMENT, LLC, a Delaware limited  
 liability company; OPENAI, L.L.C., a Delaware  
 limited liability company; OPENAI GLOBAL, LLC, a  
 Delaware limited liability company; OAI  
 CORPORATION, a Delaware corporation; OPENAI  
 HOLDINGS, LLC, a Delaware limited liability  
 company; OPENAI HOLDCO, LLC, a Delaware  
 limited liability company; OPENAI INVESTMENT

Case No.: 4:22-cv-06823-JST  
 4:22-cv-07074-JST

**SECOND AMENDED COMPLAINT**

**CLASS ACTION**

**DEMAND FOR JURY TRIAL**

**\*\*FILED UNDER SEAL\*\***

**[CONFIDENTIAL]**

1 LLC, a Delaware limited liability company; OPENAI  
2 STARTUP FUND SPV I, L.P. , a Delaware limited  
partnership; and OPENAI STARTUP FUND SPV GP  
I, L.L.C. , a Delaware limited liability company;

3 *Defendants.*

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

TABLE OF CONTENTS

1

2 I. OVERVIEW: A BRAVE NEW WORLD OF SOFTWARE PIRACY ..... 1

3 II. JURISDICTION AND VENUE..... 4

4 III. INTRADISTRICT ASSIGNMENT ..... 4

5 IV. PARTIES..... 4

6 A. Plaintiffs ..... 4

7 B. Defendants ..... 6

8 V. AGENTS AND CO-CONSPIRATORS ..... 9

9 VI. CLASS ALLEGATIONS ..... 9

10 A. Class Definitions..... 9

11 B. Numerosity.....11

12 C. Typicality.....11

13 D. Commonality & Predominance.....11

14 1. DMCA Violations.....11

15 2. Contract-Related Conduct .....11

16 3. Injunctive Relief..... 12

17 4. Defenses ..... 12

18 E. Adequacy..... 12

19 F. Other Class Considerations ..... 12

20 VII. FACTUAL ALLEGATIONS ..... 13

21 A. Introduction..... 13

22 B. Codex Outputs Copyrighted Materials Without Following the Terms of the

23 Applicable Licenses ..... 13

24 C. Copilot Outputs Copyrighted Materials Without Following the Terms of the

25 Applicable Licenses ..... 17

26 D. Codex and Copilot Were Trained on Copyrighted Materials Offered Under

27 Licenses.....20

28

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

E. Copilot Was Launched Despite Its Propensity for Producing Unlawful Outputs ..... 21

F. Copilot Reproduces the Code of the Named Plaintiffs Without Attribution ..... 24

    1. Example: Copilot Outputs the Code of Doe 2 Essentially Verbatim..... 24

    2. Example: Copilot Outputs the Code of Doe 1 in Modified Format ..... 26

    3. Example: Copilot Outputs the Code of Doe 5 In Modified Format..... 30

    4. Example: Copilot Outputs Code of Doe 5 Essentially Verbatim..... 33

G. Codex and Copilot Were Designed to Withhold Attribution, Copyright Notices,  
and License Terms from Their Users ..... 37

B. Open-Source Licenses Began to Appear in the Early 1990s ..... 42

H. Microsoft Has a History of Flouting Open-Source License Requirements ..... 43

I. GitHub Was Designed to Cater to Open-Source Projects ..... 45

J. OpenAI Is Intertwined with Microsoft and GitHub..... 47

K. Conclusion of Factual Allegations ..... 49

VIII. CLAIMS FOR RELIEF ..... 49

IX. DEMAND FOR JUDGMENT ..... 58

X. JURY TRIAL DEMANDED ..... 59

1 Plaintiffs J. Doe 1, J. Doe 2, J. Doe 3, J. Doe 4 and J. Doe 5 (“Plaintiffs”), on behalf of themselves  
 2 and all others similarly situated, bring this Class Action Complaint (the “Complaint”) against Defendants  
 3 GitHub, Inc.; Microsoft Corporation; OpenAI, Inc.; OpenAI, L.P.; OpenAI OpCo, L.L.C.; OpenAI GP,  
 4 L.L.C.; OpenAI Startup Fund GP I, L.L.C.; OpenAI Startup Fund I, L.P.; OpenAI Startup Fund  
 5 Management, LLC; OpenAI, L.L.C.; OpenAI Global, LLC (“OpenAI Global”); OAI Corporation  
 6 (“OAI”); OpenAI Holdings, LLC (“OpenAI Holdings”); OpenAI Holdco, LLC; OpenAI Investment  
 7 L.L.C.; OpenAI Startup Fund SPV I, L.P.; and OpenAI Startup Fund SPV GP I, LLC<sup>1</sup> for violation of the  
 8 Digital Millennium Copyright Act, 17 U.S.C. §§ 1201–1205 (the “DMCA”); breach of contract regarding  
 9 the Suggested Licenses, and breach of contract regarding GitHub’s policies including its terms of service.

## 10 I. OVERVIEW: A BRAVE NEW WORLD 11 OF SOFTWARE PIRACY

12 1. Plaintiffs and Class members are owners of copyright interests in materials made available  
 13 publicly on GitHub that are subject to various licenses containing conditions for use of those works (the  
 14 “Licensed Materials”). All the licenses at issue here (the “Licenses”) contain certain common terms (the  
 15 “License Terms”).

16 2. “Artificial Intelligence” is referred to herein as “AI.” AI is defined for the purposes of this  
 17 Complaint as a computer program that algorithmically simulates human reasoning or inference, often  
 18 using statistical methods. Machine Learning (“ML”) is a subset of AI in which the behavior of the  
 19 program is derived from analyzing a corpus of material called training data.

20 3. GitHub is a company founded in 2008 by a team of open-source enthusiasts. At the time,  
 21 GitHub’s stated goal was to support open-source development, especially by hosting open-source source

---

22  
 23 <sup>1</sup> GitHub, Inc. is referred to as “GitHub.” Microsoft Corporation is referred to as “Microsoft.” OpenAI,  
 24 Inc.; OpenAI, L.P.; OpenAI OpCo, L.L.C.; OpenAI GP, L.L.C.; OpenAI Startup Fund GP I, L.L.C.;  
 25 OpenAI Startup Fund I, L.P.; OpenAI Startup Fund Management, LLC; OpenAI, L.L.C.; OpenAI  
 26 Global; OAI Corporation (“OAI”); OpenAI Holdings; OpenAI Holdco, LLC; OpenAI Investment LLC;  
 27 OpenAI Startup Fund SPV I, L.P.; and OpenAI Startup Fund SPV GP I, L.L.C. are referred to  
 28 collectively herein as “OpenAI.” Collectively, GitHub, Inc., Microsoft Corporation, OpenAI, Inc.;  
 OpenAI, L.P.; OpenAI GP, L.L.C.; OpenAI Startup Fund GP I, L.L.C.; OpenAI Startup Fund I, L.P.;  
 OpenAI Startup Fund Management, LLC; OpenAI, L.L.C.; OpenAI Global; OAI; OpenAI Holdings;  
 OpenAI Holdco, LLC; OpenAI Investment LLC; OpenAI Startup Fund SPV I, L.P.; and OpenAI Startup  
 Fund SPV GP I, L.L.C. are referred to herein as “Defendants.”

1 code on the website github.com. Over the next 10 years, GitHub, based on these representations  
2 succeeded wildly, attracting nearly 25 million developers.

3 4. Developers published Licensed Materials on GitHub pursuant to written Licenses. In  
4 particular, the most popular ones share a common term: use of the Licensed Materials requires some form  
5 of *attribution*, usually by, among other things, including a copy of the license along with the name and  
6 copyright notice of the original author.

7 5. On October 26, 2018, Microsoft acquired GitHub for \$7.5 billion. Though some members  
8 of the open-source community were skeptical of this union, Microsoft repeated one mantra throughout:  
9 “Microsoft Loves Open Source.” For the first few years, Microsoft’s representations seemed credible.

10 6. Microsoft invested \$1 billion in OpenAI LP in July 2019 at a \$20 billion valuation. In 2020,  
11 Microsoft became exclusive licensee of OpenAI’s GPT-3 language model—despite OpenAI’s continued  
12 claims its products are meant to benefit “humanity” at large. In 2021, Microsoft began offering GPT-3  
13 through its Azure cloud-computing platform. On October 20, 2022, it was reported that OpenAI “is in  
14 advanced talks to raise more funding from Microsoft” at that same \$20 billion valuation. Copilot runs on  
15 Microsoft’s Azure platform. Microsoft has used Copilot to promote Azure’s processing power,  
16 particularly regarding AI.

17 7. On information and belief, Microsoft obtained a partial ownership interest in OpenAI in  
18 exchange for its \$1 billion investment. As OpenAI’s largest investor and largest service provider—  
19 specifically in connection with Microsoft’s Azure product—Microsoft exerts considerable control over  
20 OpenAI.

21 8. In June 2021, GitHub and OpenAI launched Copilot, an AI-based product that promises to  
22 assist software coders by providing or filling in blocks of code using AI. GitHub charges Copilot users \$10  
23 per month or \$100 per year for this service. Copilot ignores, violates, and removes the Licenses offered by  
24 thousands—possibly millions—of software developers, thereby accomplishing software piracy on an  
25 unprecedented scale. Copilot outputs text derived from Plaintiffs’ and the Class’s Licensed Materials  
26 without adhering to the applicable License Terms and applicable laws. Copilot’s output is referred herein  
27 as “Output.”

1           9.       On August 10, 2021, OpenAI debuted its Codex product, which converts natural language  
2 into code and is integrated into Copilot. Copilot and Codex can be called either AIs or MLs. Codex and  
3 Copilot will be referred to as AIs herein unless a distinction is required.

4           10.      Though Defendants have been cagey about what data was used to train the AI,<sup>2</sup> they have  
5 conceded that the training data includes data in vast numbers of publicly accessible repositories on  
6 GitHub,<sup>3</sup> which include and are limited by Licenses.

7           11.      Among other things, Defendants stripped Plaintiffs' and the Class's attribution, copyright  
8 notice, and license terms from their code in violation of the Licenses and Plaintiffs' and the Class's rights.  
9 Defendants used Copilot to distribute the now-anonymized code to Copilot users as if it were created by  
10 Copilot.

11          12.      Copilot is run entirely on Microsoft's Azure cloud-computing platform.

12          13.      Copilot often simply reproduces code that can be traced back to open-source repositories  
13 or open-source licensees. Contrary to and in violation of the Licenses, code reproduced by Copilot *never*  
14 includes attributions to the underlying authors.

15          14.      GitHub and OpenAI have offered shifting accounts of the source and amount of the code  
16 or other data used to train and operate Copilot. They have also offered shifting justifications for why a  
17 commercial AI product like Copilot should be exempt from these license requirements, often citing "fair  
18 use."

19          15.      It is not fair, permitted, or justified. On the contrary, Copilot's goal is to replace a huge  
20 swath of open source by taking it and keeping it inside a GitHub-controlled paywall. It violates the licenses  
21 that open-source programmers chose and monetizes their code despite GitHub's pledge never to do so.

---

25 <sup>2</sup> "Training" an AI, as described in greater detail below, means feeding it large amounts of data that it  
26 interprets using given criteria. Feedback is then given to it to fine-tune its Output until it can provide  
27 Output with minimal errors.

28 <sup>3</sup> Repositories are containers for individual coding projects. They are where GitHub users upload their  
code and where other users can find it. Most GitHub users have multiple repositories.

[CONFIDENTIAL]

## II. JURISDICTION AND VENUE

1  
2 16. Plaintiffs bring this action on their own behalf as well as representatives of a Class of  
3 similarly situated individuals and entities. They seek to recover injunctive relief and damages as a result  
4 and consequence of Defendants' unlawful conduct.

5 17. Jurisdiction and venue are proper in this judicial district under 28 U.S.C. § 1331 pursuant  
6 to Defendants' violation of Section 1202(b) of the Digital Millennium Copyright Act, 17 U.S.C. §§ 1201-  
7 1205; and because a substantial part of the events giving rise to Plaintiffs' claims occurred in this District,  
8 a substantial portion of the affected interstate trade and commerce was carried out in this District, and  
9 three or more of the Defendants reside in this District and/or are licensed to do business in this District.  
10 Each Defendant has transacted business, maintained substantial contacts, and/or committed overt acts in  
11 furtherance of the illegal scheme and conspiracy throughout the United States, including in this District.  
12 Defendants' conduct has had the intended and foreseeable effect of causing injury to persons residing in,  
13 located in, or doing business throughout the United States, including in this District.

## III. INTRADISTRICT ASSIGNMENT

14  
15 18. Pursuant to Civil Local Rule 3.2 (c) and (e), assignment of this case to the San Francisco  
16 Division of the United States District Court for the Northern District of California is proper because a  
17 substantial amount of the development of the Copilot product as well as of the interstate trade and  
18 commerce involved and affected by Defendants' conduct giving rise to the claims herein occurred in this  
19 Division. Furthermore, Defendants GitHub and all the OpenAI entities are headquartered within this  
20 Division.  
21

## IV. PARTIES

### A. Plaintiffs

22  
23  
24 19. Plaintiff J. Doe 1, [REDACTED], is a resident of the State of New Hampshire. Plaintiff Doe 1  
25 published Licensed Materials they owned a copyright interest in to at least one GitHub repository under  
26 one of the Suggested Licenses. Specifically, Doe 1 has published Licensed Materials they claim a  
27 copyright interest in under the following Suggested Licenses: MIT License and GNU General Public  
28



[CONFIDENTIAL]

1 License version 3.0. Plaintiff was, and continues to be, injured during the Class Period as a result of  
2 Defendants' unlawful conduct alleged herein.

3 20. Plaintiff J. Doe 2, [REDACTED], is a resident of the State of Illinois. Plaintiff Doe 2 published  
4 Licensed Materials they owned a copyright interest in to at least one GitHub repository under one of the  
5 Suggested Licenses. Specifically, Doe 2 has published Licensed Materials they claim a copyright interest  
6 in under the following Suggested Licenses: MIT License; GNU General Public License version 3.0; GNU  
7 Affero General Public License version 3.0; The 3-Clause BSD License; and Apache License 2.0. Plaintiff  
8 was, and continues to be, injured during the Class Period as a result of Defendants' unlawful conduct  
9 alleged herein.

10 21. Plaintiff J. Doe 3, [REDACTED], is a resident of the State of Idaho. Plaintiff Doe 3  
11 published Licensed Materials they owned a copyright interest in to at least one GitHub repository under  
12 one of the Suggested Licenses. Specifically, Doe 3 has published Licensed Materials they claim a  
13 copyright interest in under the following Suggested Licenses: MIT License; GNU General Public License  
14 version 3.0; and GNU Affero General Public License version 3.0. Plaintiff was, and continues to be,  
15 injured during the Class Period as a result of Defendants' unlawful conduct alleged herein.

16 22. Plaintiff J. Doe 4, [REDACTED], is a resident of the State of South Carolina. Plaintiff Doe 4  
17 published Licensed Materials they owned a copyright interest in to at least one GitHub repository under  
18 one of the Suggested Licenses. Specifically, Doe 4 has published Licensed Materials they claim a  
19 copyright interest in under the following Suggested Licenses: GNU General Public License v2.0 and  
20 GNU General Public License v3.0. Plaintiff was, and continues to be, injured during the Class Period as a  
21 result of Defendants' unlawful conduct alleged herein.

22 23. Plaintiff J. Doe 5, [REDACTED], is a resident of the Commonwealth of Massachusetts.  
23 Plaintiff Doe 5 published Licensed Materials they owned a copyright interest in to at least one GitHub  
24 repository under one of the Suggested Licenses. Specifically, Doe 5 has published Licensed Materials they  
25 claim a copyright interest in under the following Suggested Licenses: MIT License; Apache License 2.0;  
26 and GNU General Public License v3.0.

**B. Defendants**

24. Defendant GitHub, Inc. is a Delaware corporation with its principal place of business located at 88 Colin P Kelly Jr Street, San Francisco, CA 94107. GitHub sells, markets, and distributes Copilot throughout the internet and other sales channels throughout the United States, including in this District. GitHub released Copilot on a limited “technical preview” basis on June 29, 2021. On June 21, 2022, Copilot was released to the public as a subscription-based service for individual developers. GitHub is a party to the unlawful conduct alleged herein.

25. Defendant Microsoft Corporation is a Washington corporation with its principal place of business located at One Microsoft Way, Redmond, Washington 98052. Microsoft announced its acquisition of Defendant GitHub, Inc. on June 4, 2018. On October 26, 2018, Microsoft finalized its acquisition of GitHub. Microsoft owns and operates GitHub. Through its corporate ownership, control of the GitHub Board of Directors, active management, and other means, Microsoft sells, markets, and distributes Copilot. Microsoft is a party to the unlawful conduct alleged herein.

26. Defendant OpenAI, Inc. is a Delaware nonprofit corporation with its principal place of business located at 3180 18th Street, San Francisco, CA 94110. OpenAI, Inc. is a party to the unlawful conduct alleged herein. It—along with OpenAI, L.P.—programed, trained, and maintains Codex, which infringes all the same rights as Copilot and is also an integral piece of Copilot. Copilot requires Codex to function. OpenAI, Inc. is a party to the unlawful conduct alleged herein. OpenAI, Inc. founded, owns, and exercises control over all the other OpenAI entities, including those set forth in Paragraphs 27–40.

27. 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 party to the unlawful conduct alleged herein. Its primary activity is research and technology. OpenAI, L.P. is a wholly owned subsidiary of OpenAI, Inc. that is operated for profit. OpenAI, L.P. is the OpenAI entity that co-created Copilot and offers it jointly with GitHub. OpenAI’s revenue, including revenue from Copilot, is received by OpenAI, L.P. OpenAI, Inc. controls OpenAI, L.P. directly and through the other OpenAI entities.

28. 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 party to the unlawful conduct alleged herein. Its primary activity is research and technology. OpenAI OpCo,

1 L.L.C. is a wholly owned subsidiary of OpenAI, Inc. that is operated for profit. OpenAI OpCo, L.L.C. is  
2 the OpenAI entity that co-created Copilot and offers it jointly with GitHub. OpenAI’s revenue, including  
3 revenue from Copilot, is received by OpenAI OpCo, L.L.C. OpenAI, Inc. controls OpenAI OpCo, L.L.C.  
4 directly and through the other OpenAI entities.

5 29. Defendant OpenAI GP, L.L.C. (“OpenAI GP”) is a Delaware limited liability company  
6 with its principal place of business located at 3180 18th Street, San Francisco, CA 94110. OpenAI GP is  
7 the general partner of OpenAI, L.P. OpenAI GP manages and operates the day-to-day business and affairs  
8 of OpenAI, L.P. OpenAI GP is liable for the debts, liabilities and obligations of OpenAI, L.P., including  
9 litigation and judgments. OpenAI GP is a party to the unlawful conduct alleged herein. Its primary activity  
10 is research and technology. OpenAI GP is the general partner of OpenAI, L.P. OpenAI GP was aware of  
11 the unlawful conduct alleged herein and exercised control over OpenAI, L.P. throughout the Class Period.  
12 OpenAI, Inc. directly controls OpenAI GP. OpenAI GP directly controls OpenAI Holdings and OpenAI  
13 Global.

14 30. Defendant OpenAI Startup Fund I, L.P. (“OpenAI Startup Fund I”) is a Delaware limited  
15 partnership with its principal place of business located at 3180 18th Street, San Francisco, CA 94110.  
16 OpenAI Startup Fund I was instrumental in the foundation of OpenAI, L.P., including the creation of its  
17 business strategy and providing initial funding. Through participation in OpenAI Startup Fund I, certain  
18 entities and individuals obtained an ownership interest in OpenAI, L.P. Plaintiffs are informed and  
19 believed, and on that basis allege that OpenAI Startup Fund I participated in the organization and  
20 operation of OpenAI, L.P. OpenAI Startup Fund I is a party to the unlawful conduct alleged herein.  
21 OpenAI Startup Fund I was aware of the unlawful conduct alleged herein and exercised control over  
22 OpenAI, L.P. throughout the Class Period.

23 31. Defendant OpenAI Startup Fund GP I, L.L.C. (“OpenAI Startup Fund GP I”) is a  
24 Delaware limited liability company with its principal place of business located at 3180 18th Street, San  
25 Francisco, CA 94110. OpenAI Startup Fund GP I is the general partner of OpenAI Startup Fund I.  
26 OpenAI Startup Fund GP I manages and operates the day-to-day business and affairs of OpenAI Startup  
27 Fund I. OpenAI Startup Fund GP I is liable for the debts, liabilities and obligations of OpenAI Startup  
28 Fund I, including litigation and judgments. OpenAI Startup Fund GP I was aware of the unlawful conduct

1 alleged herein and exercised control over OpenAI, L.P. throughout the Class Period. OpenAI Startup  
2 Fund GP I is a party to the unlawful conduct alleged herein. Sam Altman, co-founder, CEO, and Board  
3 member of OpenAI, Inc. is the Manager of OpenAI Startup Fund GP I. OpenAI Startup Fund GP I is the  
4 General Partner of OpenAI Startup Fund I, L.P.

5 32. Defendant OpenAI Startup Fund Management, LLC (“OpenAI Startup Fund  
6 Management”) is a Delaware limited liability company with its principal place of business located at 3180  
7 18th Street, San Francisco, CA 94110. OpenAI Startup Fund Management is a party to the unlawful  
8 conduct alleged herein. OpenAI Startup Fund Management was aware of the unlawful conduct alleged  
9 herein and exercised control over OpenAI, L.P. throughout the Class Period.

10 33. Defendant OpenAI, L.L.C. is a Delaware limited liability company with its principal place  
11 of business in San Francisco, California. OpenAI LLC owns some or all of the services and products  
12 provided by OpenAI. The sole member of OpenAI, L.L.C. is Defendant OpenAI OpCo, L.L.C.

13 34. Defendant OpenAI Global, LLC is a Delaware limited liability company with its principal  
14 place of business in San Francisco, California. OpenAI Global’s only members are Microsoft and  
15 Defendant OAI Corporation. Microsoft owns 49% of OpenAI Global, and exercises control over it as its  
16 largest minority shareholder. OpenAI describes OpenAI Global as a “capped profit company”),

17 35. Defendant OAI Corporation (“OAI”) is a Delaware corporation with its principal place of  
18 business in San Francisco, California. OAI’s only member is Defendant OpenAI Holdings LLC.

19 36. Defendant OpenAI Holdings, LLC is a Delaware limited liability company with its  
20 principal place of business in San Francisco, California. The members of OpenAI Holdings are Defendant  
21 OpenAI, Inc. and Aestas LLC, an OpenAI-related limited liability company that is not named as a  
22 defendant as of December 22, 2023. OpenAI Holdings is partially owned by OpenAI employees and  
23 outside investors.

24 37. Defendant OpenAI Holdco, LLC is a Delaware limited liability company with its principal  
25 place of business in San Francisco, California.

26 38. Defendant OpenAI Investment LLC is a Delaware limited liability company with its  
27 principal place of business in San Francisco, California.

39. Defendant OpenAI Startup Fund SPV I, L.P. is a Delaware limited partnership with its principal place of business in San Francisco, California.

40. Defendant OpenAI Startup Fund SPV GP I, L.L.C. is a Delaware limited liability company with its principal place of business in San Francisco, California. OpenAI Startup Fund SPV GP I, L.L.C. is the general partner and controls OpenAI Startup Fund SPV I, L.P.

**V. AGENTS AND CO-CONSPIRATORS**

41. The unlawful acts alleged against the Defendants in this class action complaint were authorized, ordered, or performed by the Defendants’ respective officers, agents, employees, representatives, or shareholders while actively engaged in the management, direction, or control of the Defendants’ businesses or affairs.

42. The Defendants’ agents operated under the explicit and apparent authority of their principals.

43. Each Defendant, and its subsidiaries, affiliates and agents operated as a single unified entity.

44. Various persons and/or firms not named as Defendants herein may have participated as coconspirators in the violations alleged herein and may have performed acts and made statements in furtherance thereof.

45. Each acted as the principal, agent, or joint venture of, or for other Defendants with respect to the acts, violations, and common course of conduct alleged herein.

**VI. CLASS ALLEGATIONS**

**A. Class Definitions**

46. Plaintiffs bring this action for damages and injunctive relief on behalf of themselves and all others similarly situated as a class action pursuant to Rules 23(a), 23(b)(2), and 23(b)(3) of the Federal Rules of Civil Procedure, on behalf of the following Classes:

1 **“Injunctive Relief Class” under Rule 23(b)(2):**

2 All persons or entities domiciled in the United States that, (1) owned an  
3 interest in at least one US copyright in any work; (2) offered that work  
4 under one of GitHub’s Suggested Licenses<sup>4</sup>; and (3) stored Licensed  
5 Materials in any public GitHub repositories at any time between January 1,  
6 2015 and the present (the “Class Period”).

7 **“Damages Class” under Rule 23(b)(3):**

8 All persons or entities domiciled in the United States that, (1) owned an  
9 interest in at least one US copyright in any work; (2) offered that work  
10 under one of GitHub’s Suggested Licenses; and (3) stored Licensed  
11 Materials in any public GitHub repositories at any time during the Class  
12 Period.

13 These “Class Definitions” specifically exclude the following person or entities:

- 14 a. Any of the Defendants named herein;
- 15 b. Any of the Defendants’ co-conspirators;
- 16 c. Any of Defendants’ parent companies, subsidiaries, and affiliates;
- 17 d. Any of Defendants’ officers, directors, management, employees, subsidiaries,  
18 affiliates, or agents;
- 19 e. All governmental entities; and
- 20 f. The judges and chambers staff in this case, as well as any members of their  
21 immediate families.

---

22 <sup>4</sup> When a GitHub user creates a new repository, they have the option of selecting one of thirteen licenses  
23 from a dropdown menu to apply to the contents of that repository. (They can also apply a different license  
24 later, or no license.) The Creative Commons Zero v1.0 Universal and the Unlicense donate the covered  
25 work to the public domain and/or otherwise waive all copyrights and related rights. Because they do not  
26 contain the necessary provisions nor do they even allow the owner to make copyright claims in most  
27 circumstances, they are not included in the Class Definition. We refer to the remaining eleven options as  
28 the “Suggested Licenses,” which are: (1) Apache License 2.0 (“Apache 2.0”); (2) GNU General Public  
License version 3 (“GPL-3.0”); (3) MIT License (“MIT”); (4) The 2-Clause BSD License (“BSD 2”);  
(5) The 3-Clause BSD License (“BSD 3”); (6) Boost Software License (“BSL-1.0”); (7) Eclipse Public  
License 2.0 (“EPL-2.0”); (8) GNU Affero General Public License version 3 (“AGPL-3.0”); (9) GNU  
General Public License version 2 (“GPL-2.0”); (10) GNU Lesser General Public License version 2.1  
 (“LGPL-2.1”); and (11) Mozilla Public License 2.0 (“MPL-2.0”). These Suggested Licenses each  
contain at least three common requirements for use of the Licensed Materials in a derivative work or  
copy: attribution to the owner of the Licensed Materials (“Attribution”), inclusion of a copyright notice  
 (“Copyright Notice”), and inclusion of the applicable Suggested License’s text (“License Terms”).

**B. Numerosity**

47. Plaintiffs do not know the exact number of Class members, because such information is in the exclusive control of Defendants. Plaintiffs are informed and believe that there are at least thousands of Class members geographically dispersed throughout the United States such that joinder of all Class members in the prosecution of this action is impracticable.

**C. Typicality**

48. Plaintiffs' claims are typical of the claims of their fellow Class members because Plaintiffs and Class members all own code published under a License. Plaintiffs and the Class published work subject to a License to GitHub later used by Copilot. Plaintiffs and absent Class members were damaged by this and other wrongful conduct of Defendants as alleged herein. Damages and the other relief sought herein is common to all members of the Class.

**D. Commonality & Predominance**

49. Numerous questions of law or fact common to the entire Class arise from Defendants' conduct—including, but not limited to those identified below:

**1. DMCA Violations**

- Whether Defendants' conduct violated the Class's rights under the DMCA when GitHub and OpenAI caused Codex and Copilot to ingest and distribute Licensed Materials without including any associated Attribution, Copyright Notice, or License Terms.

**2. Contract-Related Conduct**

- Whether Defendants violated the Licenses governing use of the Licensed Materials by using them to train Copilot and for republishing those materials without appending the required Attribution, Copyright Notice, or License Terms.
- Whether Defendants interfered in prospective economic relations between the Class and the public regarding the Licensed Materials by concealing the License Terms.
- Whether Defendants intentionally or negligently interfered with a prospective economic advantage.

1                   **3. Injunctive Relief**

- 2                   • Whether this Court should enjoin Defendants from engaging in the unlawful conduct  
3                   alleged herein. And what the scope of that injunction would be.

4                   **4. Defenses**

- 5                   • Whether any affirmative defense excuses Defendants’ conduct.  
6                   • Whether any statutes of limitation limit Plaintiffs’ and the Class’s potential for  
7                   recovery.  
8                   • Whether any applicable statutes of limitation should be tolled as a result of Defendants’  
9                   fraudulent concealment of their unlawful conduct.

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

12                  **E. Adequacy**

13                  51. Plaintiffs will fairly and adequately represent the interests of the Class because they have  
14 experienced the same harms as the Class and have no conflicts with any other members of the Class.  
15 Furthermore, Plaintiffs have retained sophisticated and competent counsel (“Class Counsel”) who are  
16 experienced in prosecuting Federal and state class actions throughout the United States and other  
17 complex litigation and have extensive experience advising clients and litigating intellectual property,  
18 competition, contract, and privacy matters.

19                  **F. Other Class Considerations**

20                  52. Defendants have acted on grounds generally applicable to the Class, thereby making final  
21 injunctive relief appropriate with respect to the Class as a whole.

22                  53. This class action is superior to alternatives, if any, for the fair and efficient adjudication of  
23 this controversy. Prosecuting the claims pleaded herein as a class action will eliminate the possibility of  
24 repetitive litigation. There will be no material difficulty in the management of this action as a class action.

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



## VII. FACTUAL ALLEGATIONS

### A. Introduction

55. This class action against Defendants concerns an OpenAI product called Codex and a GitHub product called Copilot.

56. OpenAI began development of Codex sometime after OpenAI was founded in December 2015 and released Codex on a limited basis in August 2021.

57. GitHub began development of Copilot sometime in 2019, released it on a limited basis in June 2021, and released it as a paid subscription service in June 2022.

58. Codex and Copilot are assistive AI-based systems offered to software programmers. These AI systems are each trained on a large corpus of publicly accessible software code and other materials, including all the Licensed Materials. Defendants describe these products as making coding easier by accepting a code “prompt” from a programmer and emitting a possible completion of that code, which is referred to herein as Output, as set forth above. Defendants claim Codex and Copilot do not retain copies of the materials they are trained on. In practice, however, the Output is often a near-identical reproduction of code from the training data.

59. Codex and Copilot are related. Copilot is a joint venture between GitHub and OpenAI. Codex is a standalone product released by OpenAI that also “powers GitHub Copilot, which [OpenAI] built and launched in partnership with GitHub.”<sup>5</sup> “GitHub Copilot uses the OpenAI Codex to suggest code and entire functions in real-time, right from your editor.”<sup>6</sup> The version of Codex used to power Copilot is modified from OpenAI’s standalone version. Defendants have kept secret the details of Codex’s modifications and its integration into or interaction with Copilot.

### B. Codex Outputs Copyrighted Materials Without Following the Terms of the Applicable Licenses

60. Below is an explanation of how Codex functions. When Codex is prompted with:  
function isEven(n) {

---

<sup>5</sup> See <https://openai.com/blog/openai-codex/>.

<sup>6</sup> See <https://github.com/features/copilot/>.

1 it assumes this is the beginning of a function written in the JavaScript language that will test whether a  
2 number is even.

3 61. Based on this assumption, Codex will then provide Output meant to complete the rest of  
4 the function. Based on the given prompt, it produced the following response:<sup>7</sup>

```
5 function isEven(n) {
6   if (n == 0)
7     return true;
8   else if (n == 1)
9     return false;
10  else if (n < 0)
11    return isEven(-n);
12  else
13    return isEven(n - 2);
14  }
15 console.log(isEven(50));
16 // → true
17 console.log(isEven(75));
18 // → false
19 console.log(isEven(-1));
20 // → ??
```

14 62. The function itself occupies the first ten lines. Six additional lines follow the function,  
15 beginning with “console.log(isEven(50))”. One possible explanation for Codex’s inclusion of these lines  
16 is to test the “isEven” function. Though not part of the function itself, the lines will confirm the function  
17 works for certain values. In this case, the code implies that “isEven(50)” should return the value “true”,  
18 and “isEven(75)” should return “false”. Those answers are correct.

19 63. The penultimate line indicates “isEven(-1)” should return “??”. This is an error, as  
20 “isEven(-1)” should return “false”.

21 64. Codex cannot and does not understand the meaning of software code or any other  
22 Licensed Materials. But in training, what became Codex was exposed to an enormous amount of existing  
23 software code (its “Training Data”) and—with input from its trainers and its own internal processes—  
24  
25

---

26 <sup>7</sup> Due to the nature of Codex, Copilot, and AI in general, Plaintiffs cannot be certain these examples  
27 would produce the same results if attempted following additional trainings of Codex and/or Copilot.  
28 However, these examples are representative of Codex and Copilot’s Output at the time just prior to the  
filing of this Complaint.

1 inferred certain statistical patterns governing the structure of code and other Licensed Materials. The  
2 finished version of Codex, once trained, is known as a “Model.”

3 65. When given a prompt, such as the initial prompt discussed above—“function isEven(n)  
4 {”—Codex identifies the most statistically likely completion, based on the examples it reviewed in  
5 training. Every instance of Output from Codex is derived from material in its Training Data. Most of its  
6 Training Data consisted of Licensed Materials.

7 66. Codex does not “write” code the way a human would, because it does not understand the  
8 meaning of code. Codex’s lack of understanding of code is evidenced when it emits extra code that is not  
9 relevant under the circumstances. Here, Codex was only prompted to produce a function called “isEven”.  
10 To produce its answer, Codex relied on Training Data that also appended the extra testing lines. Having  
11 encountered this function and the follow-up lines together frequently, Codex extrapolates they are all part  
12 of one function. A human with even a basic understanding of how JavaScript works would know the extra  
13 lines are not part of the function itself.

14 67. Beyond the superfluous and inaccurate extra lines, this “isEven” function also contains  
15 two major defects. First, it assumes the variable “n” holds an integer. It could contain some other kind of  
16 value, like a decimal number or text string, which would cause an error. Second, even if “n” does hold an  
17 integer, the function will trigger a memory error called a “stack overflow” for sufficiently large integers.  
18 For these reasons, experienced programmers would not use Codex’s Output.

19 68. Codex does not identify the owner of the copyright to this Output, nor any other—it has  
20 not been trained to provide Attribution. Nor does it include a Copyright Notice nor any License Terms  
21 attached to the Output. This is by design—Codex was not coded or trained to track or reproduce such  
22 data. The Output in the example above is taken from *Eloquent JavaScript* by Marijn Haverbeke.<sup>8</sup>

23 69. Here is the exercise from *Eloquent JavaScript*:

24 // Your code here.  
25

---

26 <sup>8</sup> <https://eloquentjavascript.net/code/#3.2>. *Eloquent JavaScript* is “Licensed under a Creative Commons  
27 [A]tribution-[N]oncommercial license. All code in this book may also be considered licensed under an  
28 MIT license.” See <https://eloquentjavascript.net/>. Thus, having also been posted on GitHub, the code  
Codex relied on meets the definition of Licensed Materials.

```

1 console.log(isEven(50));
  // → true
2 console.log(isEven(75));
  // → false
3 console.log(isEven(-1));
  // → ??
4

```

5 70. The exercise includes the “??” error. However, for Haverbeke’s purposes, this is not an  
6 error but a placeholder value for the reader to fill in. Codex—as a mere probabilistic model—fails to  
7 recognize this nuance. The inclusion of the double question marks confirms unequivocally that Codex  
8 took this code directly from a copyrighted source without following any of the attendant License Terms.

9 71. Haverbeke provides the following solution to the function discussed above:

```

10 function isEven(n) {
11   if (n == 0) return true;
12   else if (n == 1) return false;
13   else if (n < 0) return isEven(-n);
14   else return isEven(n - 2);
15 }
16 console.log(isEven(50));
17 // → true
18 console.log(isEven(75));
19 // → false
20 console.log(isEven(-1));
21 // → false
22

```

23 72. Aside from different line breaks—which are not semantically meaningful in JavaScript—  
24 this code for the function “isEven” is the same as what Codex produced. The tests are also the same,  
25 though in this case Haverbeke provides the right answer for “isEven(-1)”, which is “false”. Codex has  
26 reproduced Haverbeke’s Licensed Material almost verbatim, with the only difference being drawn from a  
27 different portion of those same Licensed Materials.

28 73. There are many copies of Haverbeke’s code stored in public repositories on GitHub, where  
programmers who are working through Haverbeke’s book store their answers.

74. The MIT license provides that “The above copyright notice and this permission notice  
shall be included in all copies or substantial portions of the Software.”<sup>9</sup> Any person taking this code

---

<sup>9</sup> See Appendix A for full text of the MIT License.

1 directly from *Eloquent JavaScript* would have direct access to these License Terms and know to follow  
2 them if incorporating the Licensed Materials into a derivative work and/or copying them. Codex does not  
3 provide these License Terms.

4 75. OpenAI Codex's Output would frequently, perhaps even constantly, contain Licensed  
5 Materials, i.e., it would have conditions associated with it through its associated license. In its 2021  
6 research paper about Codex called "Evaluating Large Language Models Trained on Code," OpenAI  
7 stated Codex's Output is "often incorrect" and can contain security vulnerabilities and other  
8 "misalignments" (meaning, departures from what the user requested).

9 76. Most open-source licenses require attribution of the author, notice of their copyright, and a  
10 copy of the license specifically to ensure that future coders can easily credit all previous authors and  
11 ensure they adhere to all applicable licenses. All the Suggested Licenses include these requirements.

12 77. Ultimately, Codex derives its value primarily from its ability to locate and output  
13 potentially useful Licensed Materials. And from its obfuscation of any rights associated with those  
14 materials.

### 15 **C. Copilot Outputs Copyrighted Materials Without Following the Terms of the Applicable** 16 **Licenses**

17 78. GitHub Copilot works in a similar way to OpenAI Codex. As mentioned above, a modified  
18 version of Codex is used as the engine that powers Copilot.

19 79. Copilot is installed by the end user as an extension to various code editors, including  
20 Microsoft's Visual Studio and VS Code. As the user types into the editor, their code is uploaded in real  
21 time to Microsoft's Azure cloud platform, where they become prompts for Copilot.

22 80. When we give Copilot the same prompt discussed above in Paragraph 61, "function  
23 isEven(n) {", it interprets the prompt as the beginning of a function written in the JavaScript language that  
24 will test whether a number is even, same as Codex.

25 81. However, the Output of Copilot in response to the prompt is different than Codex, namely:

```
26 function isEven(n) {  
27     return n % 2 === 0;  
28 }
```

1           82.     This function is much closer to what a human programmer might use as compared to  
2 Codex’s older, inaccurate offering. It handles all values and types of “n” correctly. It does not cause a  
3 stack overflow for larger values of “n” like the Codex Output.

4           83.     Copilot’s Output, like Codex’s, is derived from existing code. Namely, sample code that  
5 appears in the online book *Mastering JS*, written by Valeri Karpov.<sup>10</sup> Like Codex’s Output, Copilot’s is  
6 also based upon copyrighted educational material. *Mastering JS* is a set of educational exercises for  
7 programmers. Like *Eloquent JavaScript*, there are many copies of Karpov’s exercise stored in public  
8 repositories on GitHub. Programmers working through *Mastering JS* store their answers there.

9           84.     If Copilot is prompted with the name of a function that will test whether a number is prime  
10 (that is, a number that can only be evenly divided by 1 and itself), namely “function isPrime(n) {”, it  
11 returns:

```
12 function isPrime(n) {
13   if (n < 2) {
14     return false;
15   }
16   for (let i = 2; i < n; i++) {
17     if (n % i === 0) {
18       return false;
19     }
20   }
21   return true;
22 }
```

19           85.     Though this function will work, it contains an error often made by beginner programmers  
20 that makes it much slower than it could be. Namely, the loop in the middle, which checks possible  
21 divisors, does not need to check every divisor smaller than “n,” only the divisors smaller than the square  
22 root of “n”. As with Codex, Copilot has no understanding of how the code works. It knows that more  
23 functions called “isPrime” contain the portion that checks for all divisors smaller than “n”, so that is what  
24 it offers. It does not return what it “thinks” is best, it returns what it has seen *the most*. It is not writing, it  
25 is reproducing (i.e., copying).

26  
27  
28 <sup>10</sup> <https://masteringjs.io/tutorials/fundamentals/modulus/>.

1           86. Like the other examples above—and most of Copilot’s Output—this output is nearly a  
2 verbatim copy of copyrighted code. In this case, it is substantially similar to the “isPrime” function in the  
3 book *Think JavaScript* by Matthew X. Curinga et al,<sup>11</sup> which is:

```
4 function isPrime(n) {
5   if (n < 2) {
6     return false;
7   }
8   for (let i = 2; i < n; i++) {
9     if (n % i === 0) {
10      return false;
11    }
12  }
13  return true;
14 }
```

11           87. As with the other examples above, the source of Copilot’s Output is a programming  
12 textbook. Also like the books the other examples were taken from, there are many copies of Curinga’s  
13 code stored in public repositories on GitHub where programmers who are working through Curinga’s  
14 book keep copies of their answers.

15           88. The material in Curinga’s book is made available under the GNU Free Documentation  
16 License. Although this is not one of the Suggested Licenses, it contains similar attribution provisions,  
17 namely that “You may copy and distribute the Document in any medium, either commercially or  
18 noncommercially, provided that this License, the copyright notices, and the license notice saying this  
19 License applies to the Document are reproduced in all copies, and that you add no other conditions  
20 whatsoever to those of this License.”<sup>12</sup>

21           89. As with Codex, Copilot does not provide the end user any attribution of the original author  
22 of the code, nor anything about their license requirements. There is no way for the Copilot user to know  
23 that they must provide attribution, copyright notice, nor a copy of the license’s text. And with regard to  
24 the GNU Free Documentation License, Copilot users would not be aware that they are limited in what  
25 conditions they can place on the use of derivative works they make using this copyrighted code. Had the  
26

---

27 <sup>11</sup> <https://matt.curinga.com/think-js/#solving-problems-with-for-loops>.

28 <sup>12</sup> <https://matt.curinga.com/think-js/#gnu-free-documentation-license>.

1 Copilot user found this code in a public GitHub repository or a copy of the book it was originally  
2 published in, they would find the GNU Free Documentation License at the same time and be aware of its  
3 terms. Copilot finds that code for the user but excises the license terms, copyright notice, and attribution.  
4 This practice allows its users to assume that the code can be used without restriction. It cannot.

#### 5 **D. Codex and Copilot Were Trained on Copyrighted Materials Offered Under Licenses**

6 90. Codex is an AI system. Another way to describe it is a “model.” Without Codex, Copilot,  
7 or another AI-code-lookup-tool, code is written both by originating code from the writer’s own knowledge  
8 of how to write code as well as by finding pre-written portions of code that—under the terms of the  
9 applicable license—may be incorporated into the coding project.

10 91. Unlike a human programmer that has learned how code works and notices when code it is  
11 copying has attached license terms, a copyright notice, and/or attribution, Codex and Copilot were  
12 developed by feeding a corpus of material, called “training data,” into them. These AI programs ingest all  
13 the data and, through a complex probabilistic process, predict what the most likely solution to a given  
14 prompt a user would input is. Though more complicated in practice, essentially Copilot returns the  
15 solution it has found in the most projects when those projects are somehow weighted to adjust for  
16 whatever variables Codex or Copilot have identified as relevant.

17 92. Codex and Copilot were not programmed to treat attribution, copyright notices, and  
18 license terms as legally essential. Defendants made a deliberate choice to expedite the release of Copilot  
19 rather than ensure it would not provide unlawful Output.

20 93. The words “study” and “training” and “learning” in connection with AI describe  
21 algorithmic processes that are not analogous to human reasoning. AI models cannot “learn” as humans  
22 do, nor can it “understand” semantics and context the way humans do. Rather, it detects statistically  
23 significant patterns in its training data and provides Output derived from its training data when  
24 statistically appropriate. A “brute force” approach like this would not be efficient nor even possible for  
25 humans. A human could not memorize, statistically analyze, and easily access thousands of gigabytes of  
26 existing code, a task now possible for powerful computers like those that make up Microsoft’s Azure cloud  
27 platform. To accomplish the same task, a human may search for Licensed Materials that serve their  
28 purpose if they believe such materials exist. And if that human finds such materials, they will probably



1 abide by its License Terms rather than risk infringing its owners' rights. At the very least, if they  
2 incorporate those Licensed Materials into their own project without following its terms they will be doing  
3 so knowingly.

#### 4 **E. Copilot Was Launched Despite Its Propensity for Producing Unlawful Outputs**

5 94. GitHub and OpenAI have not provided much detail regarding what data Codex and  
6 OpenAI were trained on. Plaintiffs know for certain from GitHub and OpenAI's statements, that both  
7 systems were trained on publicly available GitHub repositories, with Copilot having been trained on all  
8 available public GitHub repositories.

9 95. According to OpenAI, Codex was trained on "billions of lines of source code from publicly  
10 available sources, including code in public GitHub repositories." Similarly, GitHub has described<sup>13</sup>  
11 Copilot's training material as "billions of lines of public code." GitHub researcher Eddie Aftandilian  
12 confirmed in a recent podcast<sup>14</sup> that Copilot is "train[ed] on public repos on GitHub."

13 96. In a recent customer-support message, GitHub's support department clarified certain facts  
14 about training Copilot. First, GitHub said that "training for Codex (the model used by Copilot) is done by  
15 OpenAI, not GitHub." Second, in its support message, GitHub put forward a more detailed justification  
16 for its use of copyrighted code as training data:

17 Training machine learning models on publicly available data is considered  
18 fair use across the machine learning community . . . OpenAI's training of  
19 Codex is done in accordance with global copyright laws which permit the  
20 use of publicly accessible materials for computational analysis and training  
21 of machine learning models, and do not require consent of the owner of  
22 such materials. Such laws are intended to benefit society by enabling  
23 machines to learn and understand using copyrighted works, much as  
24 humans have done throughout history, and to ensure public benefit, these  
25 rights cannot generally be restricted by owners who have chosen to make  
26 their materials publicly accessible.

27 The claim that training ML models on publicly available code is widely accepted as fair use is not true.  
28 And regardless of this concept's level of acceptance in "the machine learning community," under Federal  
law, it is illegal.

---

27 <sup>13</sup> <https://github.blog/2021-06-30-github-copilot-research-recitation/>.

28 <sup>14</sup> <https://www.se-radio.net/2022/10/episode-533-eddie-aftandilian-on-github-copilot/>.

1           97.     Former GitHub CEO Nat Friedman said in June 2021—when Copilot was released to a  
2 limited number of customers—that “training ML systems on public data is fair use.”<sup>15</sup> Friedman’s  
3 statement is pure speculation; no Court has considered the question of whether “training ML systems on  
4 public data is fair use.” The Fair Use affirmative defense is only applicable to Section 501 copyright  
5 infringement. It is not a defense to violations of the DMCA, breach of contract, nor any other claim  
6 alleged herein. It cannot be used to avoid liability here. At the same time Friedman asserted “the output  
7 [of Copilot] belongs to the operator.”

8           98.     Other open-source stakeholders have made this point already. For example, in June 2021,  
9 Software Freedom Conservancy (“SFC”), a prominent open-source advocacy organization, asked  
10 Microsoft and GitHub to provide “legal references for GitHub’s public legal positions.” No references  
11 were provided by any of the Defendants.<sup>16</sup>

12           99.     Beyond the examples above, Copilot regularly Output’s verbatim copies of Licensed  
13 Materials. For example, Copilot reproduced verbatim well-known code from the game Quake III, use of  
14 which is governed by one of the Suggested Licenses—GPL-2.<sup>17</sup>

15           100.    Copilot also reproduced code that had been released under a license that allowed its use  
16 only for free games and required attribution by including a copy of the license. Copilot did not mention  
17 nor include the underlying license when providing a copy of this code as Output.<sup>18</sup>

18           101.    Texas A&M computer-science professor Tim Davis has provided numerous examples of  
19 Copilot reproducing code belonging to him without its license or attribution.<sup>19</sup>

20           102.    GitHub concedes that in ordinary use, Copilot will reproduce identical passages of code,  
21 verbatim: “Our latest internal research shows that about 1% of the time, a suggestion [Output] may  
22 contain some code snippets longer than ~150 characters that matches” code from the training data. This  
23 standard is more limited than is necessary for copyright infringement. But even using GitHub’s own

---

24 <sup>15</sup> <https://twitter.com/natfriedman/status/1409914420579344385/>.

25 <sup>16</sup> <https://sfconservancy.org/blog/2022/feb/03/github-copilot-copyleft-gpl/>.

26 <sup>17</sup> <https://twitter.com/stefankarpinski/status/1410971061181681674/>.

27 <sup>18</sup> <https://twitter.com/ChrisGr93091552/status/1539731632931803137/>.

28 <sup>19</sup> <https://twitter.com/DocSparse/status/1581461734665367554/>.

1 metric and the most conservative possible criteria, Copilot has violated the DMCA at least tens of  
2 thousands of times.

3 103. In June 2022, Copilot had 1,200,000 users. If only 1% of users have ever received Output  
4 based on Licensed Materials and only once each, Defendants have “only” breached Plaintiffs’ and the  
5 Class’s Licenses 12,000 times. However, each time Copilot outputs Licensed Materials without  
6 attribution, the copyright notice, or the License Terms it violates the DMCA three times. Thus, even  
7 using this extreme underestimate, Copilot has “only” violated the DMCA 36,000 times. Because Copilot  
8 constantly Outputs code as a user writes, and because nearly all of Copilot’s training data was Licensed  
9 Material, this number is most likely exponentially lower than the true number of breaches and DMCA  
10 violations.

11 104. Academics are continuing to study generative AI models and their behavior. Recent  
12 academic research shows that the likelihood Plaintiffs’ or class members’ code would be emitted verbatim  
13 is only increasing. For instance the study, *Quantifying Memorization Across Neural Language Models* by  
14 Nicholas Carlini et al.,<sup>20</sup> tested multiple models by feeding prefixes of prompts based on training data into  
15 each model in order to compare the performance of models of different sizes to emit output that is  
16 identical to training data. The study concluded:

17 Large language models (LMs) have been shown to memorize parts of their  
18 training data, and when prompted appropriately, they will emit the  
19 memorized training data verbatim. ... We describe three [mathematical]  
20 relationships that quantify the degree to which LMs emit memorized  
21 training data. **Memorization significantly grows as we increase (1)**  
22 **the capacity of a model**, (2) the number of times an example has been  
23 duplicated, and (3) the number of tokens of context used to prompt the  
24 model. Surprisingly, we find the situation becomes more complicated  
25 when generalizing these results across model families. On the whole, we  
26 find that memorization in LMs is more prevalent than previously believed  
27 and will likely get worse as models continues to scale.

28 (Emphasis added). Or as simply put by the study, “Bigger Models Memorize More.”

---

28 <sup>20</sup> <https://arxiv.org/pdf/2202.07646.pdf>

1           105. In other words, as generative AI models such as Copilot increase capacity and continue to  
2 scale, it becomes **more likely** that training data will become memorized and emitted verbatim, i.e., as an  
3 exact duplicate.

4           106. Given GitHub's increasing commitment to Copilot, and the scale of growth of Copilot, it  
5 follows that Copilot is more likely to emit duplicates of memorized training data as it continues to scale.

6           107. Indeed, GitHub has announced that it is adopting OpenAI's GPT-4 model for Copilot,  
7 which is a bigger and more capable language model than the Codex-derived model.

8           108. Furthermore, the Suggested Licenses impose attribution obligations not only when  
9 Licensed Materials have been used verbatim, but also when Licensed Materials have been modified or  
10 adapted. Though Output from Copilot is often a verbatim, i.e., identical copy, even more often it is a  
11 modification: for instance, a near-identical copy that contains only semantically insignificant variations of  
12 the original Licensed Materials, or a modified copy that recreates the same algorithm. Whenever Copilot  
13 outputs Licensed Materials in a manner that qualifies as a modification, the attribution requirements of  
14 the Suggested Licenses still apply. Copilot's failure to provide the attributions for outputs that are  
15 modifications of Licensed Materials represents another enormous set of license breaches.

#### 16           **F. Copilot Reproduces the Code of the Named Plaintiffs Without Attribution**

17           109. Because Copilot was trained on all available public GitHub repositories, if Licensed  
18 Materials have been posted to a GitHub public repository, Plaintiffs and the Class can be reasonably  
19 certain it was ingested by Copilot and is sometimes returned to users as Output.

20           110. Described below are some specific examples of Copilot's unlawful behavior using Licensed  
21 Materials owned by the named Plaintiffs. These examples were emitted by Copilot after prompting  
22 Copilot.

23           111. In the examples below, original code is shaded gray, prompts to Copilot are shaded orange,  
24 and outputs from Copilot are shaded light blue.

##### 25           **1. Example: Copilot Outputs the Code of Doe 2 Essentially Verbatim**

26           112. The first example demonstrates Copilot suggesting an essentially verbatim copy of code  
27 written by Doe 2.

[CONFIDENTIAL]

1 113. [REDACTED]

2 subject to the GNU General Public License v3.0. [REDACTED]

3 [REDACTED]  
4 [REDACTED] The relevant code from the  
5 original source file is shown below:

```
6 [REDACTED]  
7 [REDACTED]  
8 [REDACTED]  
9 [REDACTED]  
10 [REDACTED]  
11 [REDACTED]  
12 [REDACTED]  
13 [REDACTED]  
14 [REDACTED]  
15 [REDACTED]  
16 [REDACTED]  
17 [REDACTED]
```

18 114. When Copilot is prompted the first few lines of Doe 2's code:

```
19 [REDACTED]  
20 [REDACTED]
```

21 Copilot suggests the following:

```
22 [REDACTED]  
23 [REDACTED]  
24 [REDACTED]  
25 [REDACTED]  
26 [REDACTED]  
27 [REDACTED]  
28 [REDACTED]
```

[CONFIDENTIAL]

1 [REDACTED]  
2 [REDACTED]  
3 [REDACTED]  
4 [REDACTED]  
5 [REDACTED]

6 115. This suggestion from Copilot is identical to Doe 2’s code, except that [REDACTED]

7 [REDACTED]

8 [REDACTED]

9 These differences in the code are cosmetic and the code is functionally equivalent; otherwise, this is a  
10 verbatim copy. Doe 2’s particular arrangement and sequencing seen in his code is distinctive expression  
11 found only in one location on GitHub: [REDACTED]

12 116. Because the Copilot suggestion is a nearly verbatim reproduction of Doe 2’s unique code,  
13 it follows that Copilot copied Doe 2’s code. Copilot therefore needed to adhere to the requirements of  
14 Doe 2’s license (GNU General Public License v3.0) for that code, including providing attribution. It does  
15 not. Copilot also did not reproduce Doe 2’s license.

16 **2. Example: Copilot Outputs the Code of Doe 1 in Modified Format**

17 117. The second example demonstrates Copilot suggesting a modified copy of code written by  
18 Doe 1. To protect Doe 1’s identity, the paragraphs describing the code will be redacted.

19 118. [REDACTED] subject to the  
20 MIT License. [REDACTED]

21 [REDACTED]

22 [REDACTED]  
23 [REDACTED]  
24 [REDACTED]  
25 [REDACTED]  
26 [REDACTED]  
27 [REDACTED]  
28 [REDACTED]

[CONFIDENTIAL]

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

[REDACTED]

[CONFIDENTIAL]

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

[REDACTED]

119. When Copilot is prompted with [REDACTED]

[REDACTED]

[REDACTED]

The first suggestion from Copilot is a modification of Doe 1's code:

[REDACTED]



[CONFIDENTIAL]

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

[REDACTED]

120. [REDACTED] do not appear in any other source file on GitHub. The only way Copilot knows how to make this suggestion is because it ingested Doe 1's source file as training data. Though the Copilot suggestion is not an exact match for Doe 1's code, it is necessarily a modification based on a copy of Doe 1's code.

121. Furthermore, many distinctive expressive features of Doe 1's code have been preserved in Copilot's suggestion. For instance, Doe 1's comments in the code (in green) are reproduced almost verbatim. [REDACTED]

[CONFIDENTIAL]

1 [REDACTED]  
2 means the same thing as this Copilot-suggested code:

3 [REDACTED]  
4  
5 122. As is apparent from a cursory glance of this example, the variations between Copilot's  
6 emitted output and Doe 1's source code are cosmetic and the code is functionally equivalent; it follows  
7 that Copilot's output is a copy of Doe 1's code.

8 123. That said, Copilot also introduces mistakes into the code. For instance, [REDACTED]

9 [REDACTED]  
10 [REDACTED]  
11 124. Still, because Copilot is reproducing Doe 1's algorithm in modified format, and the  
12 obligations in Doe 1's license (the MIT License) carry with the code even if the underlying code is  
13 modified, the Copilot suggestion needs to follow the requirements of Doe 1's license for that code,  
14 including providing attribution. It does not. Copilot also did not reproduce Doe 1's license.

15 **3. Example: Copilot Outputs the Code of Doe 5 In Modified Format**

16 125. The third example demonstrates Copilot suggesting multiple modified copies of code  
17 written by Doe 5 in response to a sequence of prompts, which is a common way of using Copilot. To  
18 protect Doe 5's identity, the paragraphs describing the code will be redacted.

19 126. [REDACTED]

20 [REDACTED] subject to the MIT License. [REDACTED]

21 [REDACTED] The relevant code from the original source file is shown below:

22 [REDACTED]  
23 [REDACTED]  
24 [REDACTED]  
25 [REDACTED]  
26 [REDACTED]  
27 [REDACTED]  
28 [REDACTED]

[CONFIDENTIAL]

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

[REDACTED]

127. When Copilot is prompted the first section of Doe 5's code, comprising the first complete test and the name of the second:

[REDACTED]

128. The first suggestion from Copilot offers to complete the prompt with a verbatim copy of Doe 5's original code, except that [REDACTED] (a variation that does not affect how the code works):

[REDACTED]

129. Next, if the name of the third test is appended, the next prompt to Copilot looks like this:

[REDACTED]

[CONFIDENTIAL]

1 [REDACTED]  
2 [REDACTED]  
3 [REDACTED]  
4 [REDACTED]  
5 [REDACTED]

6 130. The first suggestion from Copilot offers to complete the prompt with a functionally  
7 identical copy of Doe 5's code, except [REDACTED]  
8 [REDACTED] (neither of these variations affect how the code works):

9 [REDACTED]  
10 [REDACTED]  
11 [REDACTED]  
12 [REDACTED]  
13 [REDACTED]  
14 [REDACTED]  
15 [REDACTED]  
16 [REDACTED]  
17 [REDACTED]  
18 [REDACTED]  
19 [REDACTED]

20 131. As is apparent from the high degree of similarity and minor cosmetic deviations between  
21 Copilot's emitted output and Doe 5's source code, Copilot ingested, copied and reproduced Doe 5's  
22 source code as output.

23 132. Because Copilot is (repeatedly) reproducing Doe 5's original code in modified format, and  
24 the obligations in Doe 5's license (the MIT License) carries with the code even when it is modified, the  
25 Copilot suggestions need to follow the requirements of Doe 5's license for that code, including providing  
26 attribution. They do not. Copilot also did not reproduce Doe 5's license.

[CONFIDENTIAL]

**4. Example: Copilot Outputs Code of Doe 5 Essentially Verbatim**

133. The fourth example also demonstrates Copilot suggesting multiple modified copies of code written by Doe 5 in response to a sequence of prompts, which is a common way of using Copilot. To protect Doe 5's identity, the paragraphs describing the code will be redacted.

134. [REDACTED]

[REDACTED] subject to the MIT License. [REDACTED]

[REDACTED] The first three tests from the original source file are shown below:

```
[REDACTED]
```

[CONFIDENTIAL]

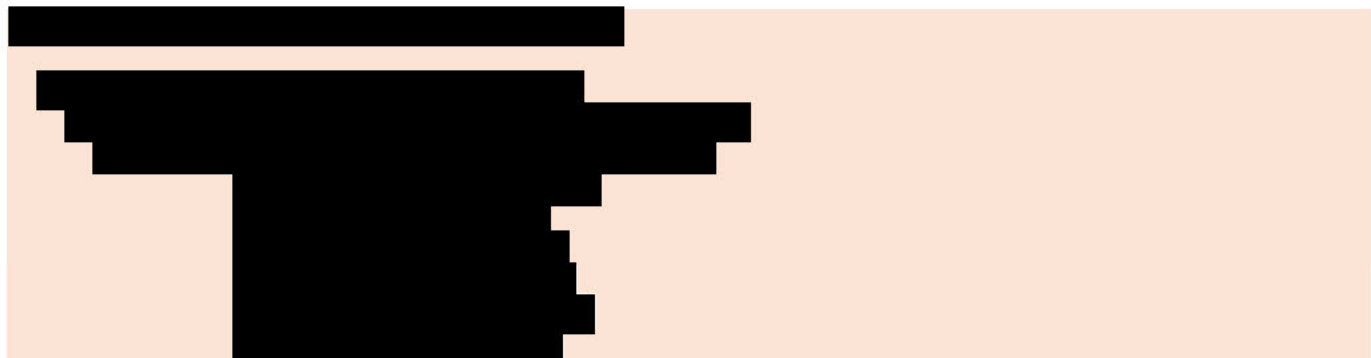
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28



135. When Copilot is prompted with the first section of Doe 5's code, comprising the first complete test and the name of the second:



The first suggestion from Copilot offers to complete the second test with a verbatim copy of Doe 5's original code:



[CONFIDENTIAL]

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

[REDACTED]

136. When Copilot's suggestion is accepted and the name of Doe 5's third test is appended, the next prompt to Copilot looks like this:

[REDACTED]

[CONFIDENTIAL]

1 [REDACTED]

2 137. Once again, the first suggestion from Copilot offers to complete the third test with a  
3 verbatim copy of Doe 5's code (except for small cosmetic variations in line breaks):

4 [REDACTED]

5 [REDACTED]  
6 [REDACTED]  
7 [REDACTED]  
8 [REDACTED]  
9 [REDACTED]  
10 [REDACTED]

11 [REDACTED]  
12 [REDACTED]

13 [REDACTED]  
14 [REDACTED]  
15 [REDACTED]  
16 [REDACTED]  
17 [REDACTED]  
18 [REDACTED]

19 [REDACTED]

20 [REDACTED]  
21 [REDACTED]  
22 [REDACTED]  
23 [REDACTED]  
24 [REDACTED]  
25 [REDACTED]

26 [REDACTED]  
27 [REDACTED]  
28 [REDACTED]



[CONFIDENTIAL]

1 [REDACTED]  
2 [REDACTED]  
3 138. Because Copilot is (repeatedly) reproducing Doe 5's code essentially verbatim, the Copilot  
4 suggestions need to follow the requirements of Doe 5's license (the MIT License) for that code, including  
5 providing attribution. They do not. Copilot also did not reproduce Doe 5's license.

6 139. These are only a few examples of Plaintiffs' code being reproduced by Copilot. It follows  
7 that many if not all prompts entered into Copilot will readily cause it to emit verbatim, near-verbatim or  
8 modified copies of Licensed Material that violate the licenses under which the source code is published.  
9 Multiplied across the many users of Copilot and the many times Copilot is prompted, each day these  
10 violations must be accruing with astonishing frequency. It is therefore likely if not certain that verbatim,  
11 near-verbatim or modified copies of each Plaintiffs' code have already been emitted by Copilot.

12 140. Additionally, even though Plaintiffs have been able to generate these examples, Plaintiffs  
13 remain at a great evidentiary disadvantage relative to Defendants, because Defendants control all the  
14 information about the training dataset. In particular, only Defendants know *when* the Licensed Materials  
15 of Plaintiffs and the Class were scraped. As is typical in open source, many of the Licensed Materials are  
16 regularly updated. As such, it is difficult to determine which iterations of code may have been trained on  
17 and would be subject to emission by Copilot.

18 **G. Codex and Copilot Were Designed to Withhold Attribution, Copyright Notices, and License**  
19 **Terms from Their Users**

20 141. Codex and Copilot have no way to determine whether license text or other Copyright  
21 Management Information ("CMI")<sup>21</sup> is part of the code it appears immediately before or after. Unless  
22 instructed otherwise, it will assume that CMI that usually appears just before a given block of code is an  
23 important part of that code or otherwise necessary for it to function.

24 142. It is a common practice to provide the applicable license text at the top of every source file  
25 in the codebase. The purpose of this practice is to avoid the code from being divorced from the license.  
26 This may occur via "vendoring," a method of creating a derivative work by including source files from a  
27

28 <sup>21</sup> CMI is defined in detail below in Paragraph 211.

1 copyrighted project directly into another project without following the terms of the license or providing  
2 attribution or a copyright notice. Copilot circumvents this protective measure to mask the degree of  
3 vendoring it engages in.

4 143. Early iterations of Copilot reproduced license text. For example, in a blog post, GitHub  
5 noted “In one instance, GitHub Copilot suggested starting an empty file with something it had even seen  
6 more than a whopping 700,000 different times during training—that was the GNU General Public  
7 License.”<sup>22</sup> Copilot no longer suggests licenses in this way because it has been altered not to. As GitHub  
8 explains: “GitHub Copilot *has* changed to require a minimum file content. So some of the suggestions  
9 flagged here would not have been shown by the current version.”

10 144. In July 2021, near Copilot’s launch, it would sometimes produce license text, attribution,  
11 and copyright notices. This CMI was not always accurate. Copilot no longer reproduces these types of  
12 CMI, incorrect or otherwise, on a regular basis. It has been altered not to.

13 145. In July 2022, in response to public criticism of Copilot’s mishandling of Licensed  
14 Materials, GitHub introduced a user-settable Copilot filter called “Suggestions matching public code.” If  
15 set to “block,” this filter claims to prevent Copilot from suggesting verbatim excerpts of “about 150  
16 characters” that come from Licensed Materials. But even assuming the filter works as advertised, because  
17 it only checks for verbatim excerpts, it does nothing to impede the Outputs from Copilot that are  
18 modifications of Licensed Materials. Thus, as a device for respecting the rights of Plaintiffs and the Class,  
19 it is essentially worthless.

20 146. GitHub Copilot now “includes an option to either allow or block code completion  
21 suggestions that match publicly available code. If you choose to block suggestions matching public code,  
22 GitHub Copilot checks code completion suggestions with their surrounding code of about 150 characters  
23 against public code on GitHub. If there is a match, or a near match, the suggestion is not shown to you.”<sup>23</sup>  
24 There is no reason provided for GitHub’s choice of 150 matching characters of code. Nor is there a reason  
25 for GitHub to include an option to block suggestions that “match” code unless Copilot is capable of (and

---

26 <sup>22</sup> <https://github.blog/2021-06-30-github-copilot-research-recitation/>.

27 <sup>23</sup> <https://docs.github.com/en/copilot/configuring-github-copilot/configuring-github-copilot-settings-on-githubcom#enabling-or-disabling-duplication-detection>

1 does) emit verbatim copies of code. Nonetheless, GitHub provides the choice to Copilot’s paying users to  
2 use code Copilot outputs that is identical to code on its public repositories, subject to open source  
3 licenses.

4 147. GitHub now admits<sup>24</sup> that “You can opt to allow GitHub Copilot to suggest code  
5 completions that match publicly available code on GitHub.com.” GitHub states that “If you have allowed  
6 suggestions that match public code, GitHub Copilot can provide you with details about the matching code  
7 when you accept such suggestions. This feature is called code referencing.” Again, the only reason  
8 GitHub would inform users that they can opt in to allowing Copilot to suggest code completions that  
9 “match” publicly available code (i.e., code that Codex and Copilot were trained on) is that Copilot is  
10 capable of, and does, emit code suggestions that are verbatim copies of code.

11 148. GitHub states it now can provide its users, at their option, a link in to the relevant identical  
12 open-source license under which such identical code was published on GitHub. Having acknowledged  
13 there is the likelihood, if not certainty that Github will produce exact copies, and providing a tool to  
14 prevent this, GitHub makes it entirely optional to users, and provides no such optionality to licensors.  
15 Thus users who want to receive identical code from GitHub or do not want to exclude it, may do so. In so  
16 doing, GitHub facilitates and encourages users to receive identical code.

17 149. As a result, Plaintiffs are informed and believe, and on that basis allege that it is likely that  
18 their licensed code is omitted by Github in violation of the open source licenses. Further Plaintiffs are  
19 informed and believe, and on that basis allege that with respect to numerous members of the class, it is  
20 certain that some identical code has been omitted by GitbHub. Plaintiffs are informed and believe, and on  
21 that basis allege that there is a substantial risk, if not certainty, that identical code will be emitted in the  
22 future. Further, given the fact that GitHub has implemented a tool to prevent this, and have made it  
23 optional, not mandatory, to users, GitHub knows, or has reason to know that identical code will be  
24 omitted in the future.

25 150. Further, Github but makes clear it is entirely up the user to add any license or attribution to  
26 the code GitHub generated for them. “The linked web page includes details of any license identified for

---

27 <sup>24</sup> [https://docs.github.com/en/copilot/using-github-copilot/finding-public-code-that-matches-github-](https://docs.github.com/en/copilot/using-github-copilot/finding-public-code-that-matches-github-copilot-suggestions)  
28 [copilot-suggestions](https://docs.github.com/en/copilot/using-github-copilot/finding-public-code-that-matches-github-copilot-suggestions)

1 the repository where the matching code was found. Having reviewed the references, *you can decide how to*  
2 *proceed.*” This link temporarily shows up in the users’ “GitHub Copilot Log view” but “[t]he GitHub  
3 Copilot log is flushed when you close the editor.”

4 151. In addition to the fact that the tool which identifies and screens identical code is an option,  
5 *at the discretion of the user*, GitHub also acknowledges that the optional tool to prevent or exclude identical  
6 code is not 100 per cent effective. Indeed, GitHub confirms it is limited in its scope and might produce  
7 Plaintiffs or class members’ matching code, but otherwise will not detect it under certain circumstances.  
8 GitHub states “code from public repositories deleted before the index was created, may not be included in  
9 the search. For the same reason, the search may return matches to code that has been deleted or moved  
10 since the index was created.”

11 152. GitHub also admits the ability and therefore the efficacy of the code referencing tool is  
12 limited, incomplete and predictably ineffective, admits that its code referencing feature is limited to  
13 finding identical public code made from indexes of GitHub public repositories *after* training Copilot, such  
14 that “code from public repositories deleted before the index was created, may not be included in the  
15 search.” Github thereby admits that Copilot trained on code posted on GitHub can be output in identical  
16 form, without the required attribution of other compliance with license terms. This may occur which  
17 nonetheless still may not be attributable in any way, even with the best intentions of Copilot’s commercial  
18 users and even given GitHub’s “code reference” feature . Simply put, GitHub admits Copilot can and  
19 does output identical matching code of Plaintiffs and class members that its own “code referencing”  
20 feature cannot detect, even on the infrequent occasion when a user exercises the option to implement the  
21 feature.

22 153. GitHub further represents to Copilot’s paying users that “Typically, matches to public  
23 code occur in less than one percent of Copilot suggestions, so you should not expect to see code  
24 references for many of the suggestions you accept.” On information and belief, Plaintiffs are informed and  
25 believe, and allege that this number is much higher.

26 154. GitHub elsewhere represents: “If you choose to allow suggestions matching public code,  
27 and you accept a suggestion for which one or more matches were found, you can click through from an  
28

1 entry in the GitHub Copilot log to view a list of references on GitHub.”<sup>25</sup> But there is no requirement to  
2 attach the license to matching code Copilot outputs for its paying users. In other words, there are many  
3 situations in which public code is emitted without compliance with license terms.

4 155. In GitHub’s hands, the propensity for small cosmetic variations in Copilot’s Output is a  
5 feature, not a bug. GitHub does so knowing that these small cosmetic variations allow CoPilot to conceal  
6 the copying of Licensed Materials and to separate the Licensed Material from the licenses. In so doing,  
7 GitHub knowingly conceals the copying of Licensed Materials with the intent and purpose of making it  
8 difficult for Plaintiffs and members of the class to identify breaches of the licenses and to enforce their  
9 rights. These small cosmetic variations mean that GitHub can deliver to Copilot customers unlimited  
10 modified copies of Licensed Materials without ever triggering Copilot’s verbatim-code filter. AI models  
11 like Copilot often have a setting called *temperature* that specifically controls the propensity for variation in  
12 their output. On information and belief, GitHub has optimized the temperature setting of Copilot to  
13 produce small cosmetic variations of the Licensed Materials as often as possible, so that GitHub can  
14 deliver code to Copilot users that *works* the same way as verbatim code, while claiming that Copilot only  
15 produces verbatim code 1% of the time. This technique of active concealment has been effective on many  
16 occasions. Copilot is an ingenious method of software piracy and concealment.

17 156. In December 2022, GitHub launched Copilot for Business. The initial terms of service  
18 included one notable extra provision compared to ordinary Copilot: a “Defense of Third Party Claims”  
19 that read:

20 GitHub will defend you against any claim by an unaffiliated third-party that  
21 your use of GitHub Copilot misappropriated a trade secret or directly  
22 infringes a patent, copyright, trademark, or other intellectual property right  
23 of a third party, up to the greater of \$500,000.00 USD or the total amount  
24 paid to GitHub for the use of GitHub Copilot during the 12 months  
25 preceding the claim. GitHub’s defense obligations do not apply if (i) the  
26 claim is based on Code that differs from a Suggestion provided by GitHub  
Copilot, (ii) you fail to follow reasonable software development review  
practices designed to prevent the intentional or inadvertent use of Code in a  
way that may violate the intellectual property or other rights of a third party,  
or (iii) you have not enabled all filtering features available in GitHub  
Copilot.

27 <sup>25</sup> <https://docs.github.com/en/copilot/configuring-github-copilot/configuring-github-copilot-settings-on-githubcom#enabling-or-disabling-duplication-detection>  
28

1 157. If Copilot had been designed to reproduce the attribution, license terms, and copyright  
2 notices of the Licensed Materials, this kind of contractual reassurance wouldn't be necessary. With this  
3 provision (since removed), GitHub acknowledged that Copilot disrupts—possibly with legal  
4 consequences—the relationship between authors and users of open-source software.

5 **B. Open-Source Licenses Began to Appear in the Early 1990s**

6 158. In 1991, software engineer Linus Torvalds began a project to create a UNIX-like operating  
7 system that would run on common PC hardware. This project became known as Linux.

8 159. To encourage adoption of his system, and persuade other programmers to contribute, he  
9 released Linux under what was then an unusual software license called the GNU General Public License,  
10 or GPL.

11 160. The GPL is a software license. But whereas most software licenses required payment,  
12 software under the GPL is provided for free. Whereas most software licenses did not include source code,  
13 GPL software always included source code. And whereas most software licenses prohibited derivative  
14 works, the GPL not only allowed it, but encouraged it.

15 161. In certain ways, however, the GPL still operated like a traditional software license. For  
16 example, consistent with copyright law, it depended on an assertion of copyright by the software author.  
17 Even though GPL software was available at no charge, the GPL contained conditions on its users as  
18 licensees.

19 162. One license requirement was that a program derived from GPL software had to redistribute  
20 certain information about that software:

21 You may copy and distribute verbatim copies of the Program's source code  
22 as you receive it, in any medium, provided that you conspicuously and  
23 appropriately publish on each copy an appropriate copyright notice and  
24 disclaimer of warranty; keep intact all the notices that refer to this General  
25 Public License and to the absence of any warranty; and give any other  
26 recipients of the Program a copy of this General Public License along with  
27 the Program.<sup>26</sup>

28 Failure to adhere to these conditions constituted a violation of the license, triggering the possibility of

---

<sup>26</sup> <https://www.gnu.org/licenses/old-licenses/gpl-1.0.en.html>.

1 legal action. Provisions of the GPL are enforceable, and many GPL licensors have sought to enforce GPL  
2 licenses through court proceedings and other litigation.

3 163. The early years of Linux paralleled the early years of the World Wide Web. The fact that  
4 Linux was free and ran on common computer hardware made it a popular choice for web servers. Because  
5 of its contrarian GPL licensing, Linux became hugely popular. A large ecosystem of other programs and  
6 tools grew around it. This contributed to the explosive growth of the web and other network services  
7 across the rest of the 1990s.

8 164. In turn, the growth of the World Wide Web made it easier for developers in different places  
9 to collaborate on software. The GPL, and licenses like it, were a natural fit for this kind of collaborative  
10 work.

11 165. Around 1998, a new name was coined as an umbrella term for these principles of software  
12 licensing and development: *open source*.

### 13 **H. Microsoft Has a History of Flouting Open-Source License Requirements**

14 166. During the 1980s and 1990s, Microsoft was primarily a software company, focusing largely  
15 on operating systems and related applications. These included its DOS operating system and later, its  
16 Windows operating system. Windows generated billions of dollars in revenue from its sale and licensing as  
17 proprietary software for desktop computers and servers. Microsoft derived substantial income from sale  
18 of licensed products and devotes substantial resources to protecting and enforcing such licenses.

19 167. Windows is a graphical operating system. It allows users to view and store files, run  
20 software and games, play videos, and provides a way to connect to the internet.

21 168. Linux represented a competitive threat to Windows. It ran on the same hardware. It  
22 performed many of the same functions. It was free. Many programmers at the time considered Linux to be  
23 functionally superior to Windows.

24 169. Microsoft has engaged in a problematic practice known as “vaporware,” where products  
25 are announced but are in fact late, never manufactured, or canceled. Typically the company promising  
26 vaporware never has any intention of providing it. The term vaporware was coined by Microsoft in 1982 in  
27 reference to the development of its Xenix operating system.

1 170. Microsoft described its anti-Linux strategy as “FUD,” standing for fear, uncertainty, and  
2 doubt. Microsoft focused extra attention to Linux’s open-source aspects.

3 171. In 1998, a source at Microsoft leaked what became known as the “Halloween Documents”,  
4 revealing Microsoft’s thinking on how to counter the competitive threat from Linux. Among other things,  
5 the documents emphasized the importance of countering the “long term developer mindshare threat”,  
6 and concluded that to defeat open source, “[Microsoft] must target a process rather than a company.”<sup>27</sup>

7 172. In 2001, Microsoft CEO Steve Ballmer said “The way the [GPL] is written, if you use any  
8 open-source software, you must make the rest of your software open source. . . . Linux is a cancer that  
9 attaches itself in an intellectual property sense to everything it touches.”<sup>28</sup> Ballmer’s summary of GPL  
10 licensing was not accurate. In 2001, Linux was being used by corporations of every size. The growth of  
11 open source up to that point, and since, has been made possible by the open-source community’s respect  
12 for and compliance with applicable licenses.

13 173. In 2001, Microsoft was the defendant in a major software-related antitrust case, *United*  
14 *States v. Microsoft Corporation*.<sup>29</sup> In this case, the U.S. Department of Justice accused Microsoft of  
15 maintaining a software monopoly by illegally imposing technical restrictions on manufacturers of personal  
16 computers, including “tying” violations related to the Internet Explorer web browser. Judge Thomas  
17 Penfield Jackson, who presided over the antitrust trial, opined that Microsoft is “a company with an  
18 institutional disdain for both the truth and for rules of law that lesser entities must respect. It is also a  
19 company whose ‘senior management’ is not averse to offering specious testimony to support spurious  
20 defenses to claims of its wrongdoing.”<sup>30</sup>

21 174. In 2007, Microsoft admitted that it tried to influence the vote of an ISO open-standards  
22 committee by offering money to certain business partners in Sweden to vote for Microsoft’s preferred  
23 outcome.<sup>31</sup>

---

24 <sup>27</sup> <http://www.catb.org/esr/halloween/halloween1.html>.

25 <sup>28</sup> <https://lwn.net/2001/0607/a/esr-big-lie.php3>.

26 <sup>29</sup> No. Civ.A. 00-1457 TPJ.

27 <sup>30</sup> *Jackson v. Microsoft Corp.*, 135 F. Supp. 2d 38 (D.D.C. 2001).

28 <sup>31</sup> <https://learn.microsoft.com/en-us/archive/blogs/jasonmatusow/open-xml-the-vote-in-sweden/>.



1           175. After observing the rapid growth of Amazon’s original cloud computing products,  
2 Microsoft has expanded its business into cloud computing, which it has branded Microsoft Azure or  
3 simply Azure. Microsoft announced Azure to developers in 2008. It was formally released in 2010. Azure  
4 uses large-scale virtualization at Microsoft data centers and offers many hundreds of services, including  
5 infrastructure as a service (“IaaS”), platform as a service (“PaaS”), compute services, Azure Active  
6 Directory, mobile services, storage services, communication services, data management, messaging,  
7 developer services, Azure AI, blockchain, and others.

### 8           **I. GitHub Was Designed to Cater to Open-Source Projects**

9           176. By 2002, Linux had become immensely popular. But the project itself had become  
10 unwieldy and had outgrown its reliance on informal systems of managing software source code (also  
11 known as *source-control systems*). The Linux community needed something better.

12           177. Linus Torvalds set about writing a new source-control system. He named his new system  
13 Git. He released it under the GPL. It quickly became the source-control system of choice for open-source  
14 programmers.

15           178. A single software project stored in Git is called a *source repository*, commonly shortened to  
16 *repository* or just *repo*. A Git source repository would typically be stored on a networked server accessible  
17 to a group of programmers.

18           179. This became less convenient, however, when programmers were distributed among  
19 multiple locations, rather than being in a single location. A Git repository could be stored on an internet-  
20 accessible server. But setting up that server hardware and being responsible for it was inconvenient and  
21 expensive.

22           180. In 2008, a group of open-source developers in San Francisco, California founded GitHub.  
23 GitHub managed internet servers that hosted Git source repositories. With an account at GitHub, an  
24 open-source developer could easily set up a Git project accessible to collaborators anywhere in the world.  
25 From early on, GitHub’s core market was open-source developers, whom it attracted by making many of  
26 its hosting services free.

27           181. Most open-source programmers used GitHub to create “public” repositories, meaning  
28 that anyone could view them & access them. GitHub also allowed programmers and organizations to

1 create “private” repositories, which were not accessible from the public GitHub website, and required  
2 password access.

3 182. Open-source licensing was integral to GitHub. GitHub encouraged open-source developers  
4 to understand and use open-source licenses for their work. Many—though not all—public repositories on  
5 GitHub carry an open-source license. By convention, this license is stored at the top level of each  
6 repository in a file called LICENSE. GitHub’s interface also includes a button on the front pages of most  
7 repositories users can click to see details of the applicable license. A human user could easily find the  
8 license in either of these locations—as could an AI anywhere near as powerful as Codex or Copilot.

9 183. Though the GPL is one of the early open-source licenses and remains common, it is not  
10 the only open-source license. Examples of other common open-source licenses include the MIT License,  
11 the Apache License, and the Berkeley Software Distribution License (all of which are included in the  
12 Suggested Licenses).

13 184. Though these licenses differ in their wording and their details, most of them share a  
14 requirement that a copy of the license be included with any copy, derivative, or redistribution of the  
15 software, and that the author’s name and copyright notice remains intact. This is not a controversial  
16 requirement of open-source licenses—indeed, it has been an integral part of the GPL for over 30 years.

17 185. There are also many public repositories on GitHub that have no license. Though GitHub  
18 has encouraged awareness of licenses among its users, it has never imposed a default license on public  
19 repositories. A public repository without a license is subject to ordinary rules of U.S. copyright.

20 186. Open-source developers flocked to GitHub. By 2018, GitHub had become the largest and  
21 most successful Git hosting service, hosting millions of users and projects.

22 187. In October 2018, Microsoft acquired GitHub for \$7.5 billion. It was important to Microsoft  
23 that programmers use GitHub. Microsoft had developed a well-deserved poor reputation because of its  
24 documented vaporware, FUD, and other business practices, including those targeted at open-source  
25 programs and programming, and open-source licensing specifically. Microsoft made false and misleading  
26 statements and omissions to assuage such concerns, including its primary mantra intended to win over the  
27 open-source community: “Microsoft Loves Open Source.”

## J. OpenAI Is Intertwined with Microsoft and GitHub

188. OpenAI, Inc. is a nonprofit corporation founded in December 2015 by a group that included Greg Brockman, Ilya Sutskever, and other AI researchers; Elon Musk, CEO of Tesla; and Sam Altman, president of Y Combinator, a tech-startup incubator with hundreds of companies in its portfolio. Musk and Altman served as co-chairs of OpenAI, Inc. One of OpenAI, Inc.'s current board members is Reid Hoffman, founder of LinkedIn, which is now a Microsoft subsidiary. Mr. Hoffman is also a member of the Microsoft Board of Directors.

189. Less than a year later, in November 2016, OpenAI first partnered with Microsoft. It described the partnership as follows: “We’re working with Microsoft to start running most of our large-scale experiments on Azure. This will make Azure the primary cloud platform that OpenAI is using for deep learning and AI, and will let us conduct more research and share the results with the world.”

190. Initially, OpenAI, Inc. held itself out as a “non-profit artificial intelligence research company” that sought to shape AI “in the way that is most likely to benefit humanity as a whole.”

191. OpenAI, Inc. reportedly secured \$1 billion in initial funding, from sources that were largely not disclosed, but included at least most of its founders.

192. OpenAI, Inc. obtained its initial source of training data from its founders’ companies. According to reporting at the time, Musk and Altman planned to “pool[] online data from their respective companies” to serve as training data for OpenAI, Inc. projects. Musk planned to contribute data from Tesla; Altman planned to have Y Combinator companies “share their data with OpenAI.”<sup>32</sup>

193. In February 2019, Altman created OpenAI, LP, a for-profit subsidiary of the nonprofit entity OpenAI, Inc. The new OpenAI, LP entity would serve as a vessel for accepting traditional outside investment in exchange for equity and distributing profits.

194. In July 2019, OpenAI, L.P. accepted a \$1 billion investment from Microsoft. In addition to cash, Microsoft would become the exclusive licensor of certain OpenAI, LP products (including Codex, described below in Paragraph 197). Also, as part of this alliance, OpenAI, LP would use Microsoft’s cloud-computing platform, Azure, exclusively to develop and host its products. Some portion of Microsoft’s

---

<sup>32</sup> <https://www.wired.com/2015/12/elon-musks-billion-dollar-ai-plan-is-about-far-more-than-saving-the-world/>.

1 investment was paid in credits for use of Azure rather than cash. Finally, Microsoft and OpenAI agreed to  
2 “jointly build new Azure AI supercomputing technologies.”

3 195. Azure is a major growth area for Microsoft. In its most recent earnings report on October  
4 25, 2022, “Azure and other cloud services” grew by 35% from the previous quarter, more than any other  
5 product.<sup>33</sup> Azure has grown rapidly since Microsoft began its partnership with OpenAI in 2016. Its  
6 revenue grew by 50% or more every quarter from 2016 through the first three quarters of 2020.

7 196. In May 2020, Microsoft and OpenAI announced they had jointly built a supercomputer in  
8 Azure that would be used exclusively by OpenAI to train its AI models. Microsoft’s influence over and  
9 frequent collaboration with OpenAI has led some to describe Microsoft as “the unofficial owner of  
10 OpenAI.”<sup>34</sup>

11 197. One of OpenAI’s projects is GPT-3, a so-called “large language model” designed to emit  
12 naturalistic text. When researchers noticed that GPT-3 could also generate software code, they started  
13 studying whether they could make a new AI model specifically trained for this purpose. This project  
14 became known as Codex.

15 198. Sometime after July 2019, OpenAI and Microsoft began collaborating on a code-  
16 completion product for GitHub that would use Codex as its underlying model. This product became  
17 known as Copilot.

18 199. On September 28, 2022, OpenAI released an image-generation AI called DALL-E-2.  
19 Much like Copilot, DALL-E-2 removes any attribution and/or copyright notice from the images it uses to  
20 create derivative works. Like with Codex, here, OpenAI ignores the rights of the owners of copyrights to  
21 images it has ingested.

22 200. In another joint project, Microsoft and OpenAI recently launched a preview of a product  
23 called “Azure OpenAI Service.”<sup>35</sup> This service will “Leverage large-scale, generative AI models with deep  
24 understandings of language and code to enable new reasoning and comprehension capabilities for building  
25

---

26 <sup>33</sup> <https://www.microsoft.com/en-us/Investor/earnings/FY-2023-Q1/press-release-webcast/>.

27 <sup>34</sup> <https://venturebeat.com/ai/what-to-expect-from-openais-codex-api/>.

28 <sup>35</sup> <https://azure.microsoft.com/en-us/products/cognitive-services/openai-service/>.

1 cutting-edge applications. Apply these coding and language models to a variety of use cases, such as  
 2 writing assistance, code generation, and reasoning over data. Detect and mitigate harmful use with built-in  
 3 responsible AI and access enterprise-grade Azure security.”

#### 4 **K. Conclusion of Factual Allegations**

5 201. Future AI products may represent a bold and innovative step forward. GitHub Copilot and  
 6 OpenAI Codex, however, do not. Defendants should not have released these products until they could  
 7 ensure that they did not constantly violate Plaintiffs’ and the Class’s intellectual-property rights, licenses,  
 8 and other rights.

9 202. Defendants have made no attempt to comply with the open-source licenses that are  
 10 attached to much of their training data. Instead, they have pretended those licenses do not exist, and  
 11 trained Codex and Copilot to do the same. By simultaneously violating the open-source licenses of tens-  
 12 of-thousands—possibly millions—of software developers, Defendants have accomplished software piracy  
 13 on an unprecedented scale. As Microsoft’s Co-Founder Bill Gates once said regarding software piracy:  
 14 “the thing you do is theft.”<sup>36</sup>

15 203. There is no inherent limitation or constraint of AI systems that made any of this necessary.  
 16 Defendants chose to build AI systems designed to enhance their own profit at the expense of a global  
 17 open-source community that they had once sought to foster and protect. GitHub and OpenAI are profiting  
 18 at the expense of Plaintiffs’ and the Class’s rights.

### 19 **VIII. CLAIMS FOR RELIEF**

#### 20 **COUNT 1** 21 **VIOLATION OF THE DIGITAL MILLENNIUM COPYRIGHT ACT** 22 **17 U.S.C. §§ 1201–1205** **(For Injunctive Relief)** **(Against All Defendants)**

23 204. Plaintiffs and the Class hereby repeat and incorporate by reference each preceding and  
 24 succeeding paragraph as though fully set forth herein.

25 205. As described herein, Defendants have intentionally removed or altered CMI from  
 26 Plaintiffs’ code in violation of Section 1202(b)(1) of the DMCA.

27  
 28 <sup>36</sup> [https://www.digibarn.com/collections/newsletters/homebrew/V2\\_01/gatesletter.html](https://www.digibarn.com/collections/newsletters/homebrew/V2_01/gatesletter.html)

1           206. As described herein, there is a substantial risk that Defendants will distribute copies of  
2 Plaintiffs' code knowing that CMI has been removed or altered while knowing or having reasonable  
3 grounds to know that it will induce, enable, facilitate, or conceal infringement in violation of Section  
4 1202(b)(3) of the DMCA.

5           207. GitHub has admitted that about 1% of the time, a suggestion may contain code snippets  
6 longer than ~150 characters that matches code from the training data. In other words, GitHub itself has  
7 admitted that Copilot can emit identical copies of code Copilot was trained on.

8           208. GitHub has implemented features which allow users to block output blocks suggestions  
9 matching public code. GitHub would not implement such a feature unless it knows that Copilot is capable  
10 of, and does, emit output that matches code found on public repositories.

11           209. Further, given GitHub's increasing commitment to growing an AI model and the scale of  
12 Copilot's code, given that academic research suggests that a model increases the likelihood of emitting  
13 training data it has "memorized," the chance that Copilot will emit code that matches code found in the  
14 training data is only increasing as the model scales. On information and belief, if Copilot has not done so  
15 already, Copilot will emit identical copies of Class members' code.

16           210. Plaintiffs and members of the Class own the copyrights to Licensed Materials used to train  
17 Codex and Copilot. Copilot was trained on millions—possibly billions—of lines of code publicly available  
18 on GitHub. Copilot runs on Microsoft's Azure cloud platform exclusively and Microsoft had input in the  
19 creation of Copilot. Microsoft is aware that Copilot ignores License Terms and that it was trained almost  
20 exclusively on Licensed Materials.

21           211. Plaintiffs and members of the Class included the following Copyright Management  
22 Information (as defined in Section 1202(c) of the DMCA) ("CMI") in the Licensed Materials:

- 23           a. copyright notices;
- 24           b. the title and other information identifying the Licensed Materials;
- 25           c. the name of, and other identifying information about, the authors of the Licensed  
26           Materials;
- 27           d. the name of, and other identifying information about, the copyright owners of the Licensed  
28           Materials;

- e. terms and conditions for use of the Licensed Materials, specifically the Suggested Licenses; and
- f. identifying numbers or symbols referring to CMI or links to CMI.

212. Defendants did not contact Plaintiffs and the Class to obtain authority to remove or alter CMI from the Licensed Materials within the meaning of the DMCA.

213. Defendants knew that they did not contact Plaintiffs and the Class to obtain authority to remove or alter CMI from the Licensed Materials within the meaning of the DMCA.

214. As part of the scheme, Defendants did not attempt to contact Plaintiffs or Class members to obtain authority to remove or alter CMI from the Licensed Materials within the meaning of the DMCA. In fact, the removal of CMI made it difficult or impossible to contact Plaintiffs and the Class to obtain authority to remove or alter CMI from the Licensed Materials within the meaning of the DMCA. Rather, Defendants removed or altered CMI from open-source code that is owned by Class members after the code was uploaded to a GitHub repository by incorporating it into Copilot with its CMI removed.

215. Without the authority of Plaintiffs and the Class, Defendants intentionally removed or altered CMI from the Licensed Materials after they were uploaded to one or more GitHub repositories.

216. Defendants had access to but were not licensed by Plaintiffs nor the Class to train any machine learning, AI, or other pseudo-intelligent computer program, algorithm, or other functional prediction engine using the Licensed Materials.

217. Defendants had access to but were not licensed by Plaintiffs nor the Class to incorporate the Licensed Materials into Copilot.

218. Defendants had access to but were not licensed by Plaintiffs nor the Class to distribute the Licensed Materials as they do through Copilot.

219. Without the authority of Plaintiffs and the Class, Defendants distributed CMI knowing that the CMI had been removed or altered without authority of the copyright owner or the law with respect to the Licensed Materials.

220. Defendants distributed copies of the Licensed Materials knowing and intending that CMI had been removed or altered without authority of the copyright owner or the law, with respect to the Licensed Materials.

1           221. Defendants removed or altered CMI from the Licensed Materials knowing and intending  
2 that it would induce, enable, facilitate, or conceal infringement of copyright.

3           222. Without the CMI associated with the Licensed Materials, Copilot users are induced or  
4 enabled to copy the Licensed Materials. Because CMI has been removed, Copilot users do not know  
5 whether Output is owned by someone else and subject to restrictions on use. Without the CMI, copyright  
6 infringement is facilitated or concealed, because Plaintiffs and the Class are prevented from knowing or  
7 learning that the Output is based upon one or more of the Licensed Materials. Use of the Licensed  
8 Materials is not infringement when the terms of the applicable Suggested License are followed. Had the  
9 CMI not been removed, Copilot users would be aware of the Licenses and their obligations under them.  
10 The terms of the applicable Suggested License would have allowed those users to use the Licensed  
11 Materials without infringement. By withholding and concealing license information and other CMI,  
12 Defendants prevented Copilot users from making non-infringing use of the Licensed Materials. This  
13 contradicts the express wishes of Plaintiffs and the Class, which are set forth explicitly in the Suggested  
14 Licenses under which the Licensed Materials are offered.

15           223. Defendants removed or altered CMI from Licensed Materials owned by Plaintiffs and the  
16 Class while possessing reasonable grounds to know that it would induce, enable, facilitate, and/or conceal  
17 infringement of copyright in violation of Sections 1202(b)(1) and 1202(b)(3) of the DMCA.

18           224. By omitting, altering and/or concealing CMI from Copilot's Output, Defendants have  
19 reasonable grounds to know that innocent infringers are induced or enabled to copy the Licensed  
20 Materials, because CMI has been removed. Without the CMI, Defendants have reasonable grounds to  
21 know copyright infringement is facilitated or concealed, because Plaintiffs and the Class have the difficult  
22 or impossible task of proving the Licensed Materials belong to them.

23           225. The profits attributable to Defendants' violation of the DMCA include the revenue from:  
24 Copilot subscription fees, sales of or subscriptions to Defendants' Copilot-related products and/or  
25 services that are used to run Copilot, hosting Copilot on Azure, and any other of Defendants' products  
26 that contain copies of the Licensed Materials without all the original CMI. The Licensed Materials add  
27 nearly all value to the Copilot product because the purpose of Copilot is to provide code and the source of  
28 that code is the Licensed Materials. Without the Licensed Materials, Copilot would not be functional.



1           226. On information and belief, Defendants could have trained Copilot to include attribution,  
2 copyright notices, and license terms when it provides Output covered by a License.

3           227. Defendants did not request or obtain permission from Plaintiffs and the Class to use the  
4 Licensed Materials for Defendants' Copilot product.

5           228. Defendants use of the Licensed Materials does not follow the requirements of the  
6 Suggested Licenses associated with the Licensed Materials. In particular, Copilot fails to provide  
7 attribution for the creator nor the owner of the Work. Copilot fails to include the required copyright notice  
8 included in the License. Copilot fails to include the applicable Suggested License's text.

9           229. Defendants are sophisticated with respect to intellectual property matters related to open-  
10 source code. Microsoft in particular has extensive experience granting licenses, obtaining licenses, and  
11 enforcing license terms. Its most recent Annual Report states:

12                   **We protect our intellectual property investments in a variety of ways.**  
13                   **We work actively in the U.S. and internationally to ensure the**  
14                   **enforcement of copyright, trademark, trade secret, and other**  
15                   **protections that apply to our software and hardware products, services,**  
16                   **business plans, and branding.** We are a leader among technology  
17                   companies in pursuing patents and currently have a portfolio of over 69,000  
18                   U.S. and international patents issued and over 19,000 pending worldwide.  
19                   While we employ much of our internally-developed intellectual property  
20                   exclusively in our products and services, we also engage in outbound  
21                   licensing of specific patented technologies that are incorporated into  
22                   licensees' products. From time to time, we enter into broader cross-license  
23                   agreements with other technology companies covering entire groups of  
24                   patents. We may also purchase or license technology that we incorporate  
25                   into our products and services. At times, we make select intellectual  
26                   property broadly available at no or low cost to achieve a strategic objective,  
27                   such as promoting industry standards, advancing interoperability,  
28                   supporting societal and/or environmental efforts, or attracting and enabling  
                    our external development community. **Our increasing engagement with**  
                    **open source software will also cause us to license our intellectual**  
                    **property rights broadly in certain situations.**

Microsoft Corporation Annual Report, Form 10-K at 27 (July 28, 2022) (emphasis added).<sup>37</sup>

23           230. GitHub, which offers the Copilot product jointly with OpenAI, also has extensive  
24 experience with the DMCA. GitHub knows or reasonably should know that the Licensed Materials it  
25 hosts are subject to copyright. It provides the language of the Suggested Licenses to users, all of which  
26  
27

28 <sup>37</sup> <https://microsoft.gcs-web.com/static-files/07cf3c30-cfc3-4567-b20f-f4b0f0bd5087/>.

1 include copyright notices. Its 2022 Transparency Report—January to June<sup>38</sup> states: “Copyright-related  
2 takedowns (which we often refer to as DMCA takedowns) are particularly relevant to GitHub because so  
3 much of our users’ content is software code and can be eligible for copyright protection.”<sup>39</sup> In the first six  
4 months of 2022, GitHub processed 1220 DMCA takedown requests. Its DMCA Takedown Policy<sup>40</sup> notes  
5 “GitHub probably never would have existed without the DMCA.”

6 231. GitHub also knows or reasonably should know the portions of the DMCA giving rise to  
7 Plaintiffs’ claim. In its 2021 Transparency Report, “Before removing content based on alleged  
8 circumvention of copyright controls (under Section 1201 of the US DMCA or similar laws in other  
9 countries), we carefully review both the legal and technical claims, and we sponsor a Developer Defense  
10 Fund to provide developers with meaningful access to legal resources.”<sup>41</sup>

11 232. GitHub is aware that Copilot’s removal of CMI is illegal. For example, it states that  
12 “publishing or sharing tools that enable circumvention are not [permitted]”<sup>42</sup> and “Distributing tools that  
13 enable circumvention is prohibited, even if their use by developers falls under the exemption [for security  
14 research].”<sup>43</sup> GitHub has also frequently published articles discussing the DMCA, its application, and the  
15 Copyright Office’s guidance on its scope and exceptions.<sup>44</sup>

16 233. Unless Defendants are enjoined from violating the DMCA, Plaintiffs and the Class will  
17 suffer great and irreparable harm by depriving them of the right to identify and control the reproduction  
18 and/or distribution of their copyrighted works, to have the terms of their open-source licenses followed,  
19 and to pursue copyright-infringement remedies. Defendants will not be damaged if they are required to  
20 comply with the DMCA. Plaintiffs and the Class are therefore entitled to an injunction barring  
21

---

22 <sup>38</sup> <https://github.blog/2022-08-16-2022-transparency-report-january-to-june/>.

23 <sup>39</sup> <https://github.blog/2022-08-16-2022-transparency-report-january-to-june/>.

24 <sup>40</sup> <https://docs.github.com/en/site-policy/content-removal-policies/dmca-takedown-policy#what-is-the-dmca/>.

25 <sup>41</sup> <https://github.blog/2022-01-27-2021-transparency-report/>.

26 <sup>42</sup> <https://github.blog/2020-11-19-take-action-dmca-anti-circumvention-and-developer-innovation/#what-dmca-exemptions-do-not-do/>.

27 <sup>43</sup> <https://github.blog/2021-11-23-copyright-office-expands-security-research-rights/>.

28 <sup>44</sup> *See, e.g.*, Footnotes 43–46.

1 Defendants from violating the DMCA and impounding any device or product that is in the custody or  
2 control of Defendants and that the court has reasonable cause to believe was involved in a violation of the  
3 DMCA.

4 234. Defendants conspired together and acted jointly and in concert pursuant to their scheme to  
5 commit the acts that violated the DMCA alleged herein.

6 235. Defendants induced (or will induce) Copilot users to unknowingly violate the DMCA by  
7 withholding attribution, licensing, and other information as described herein.

8 **COUNT 2**  
9 **BREACH OF CONTRACT—OPEN-SOURCE LICENSE VIOLATIONS**  
10 **California Common Law**  
11 **(Against All Defendants)**

12 236. Plaintiffs and the Class hereby repeat and incorporate by reference each preceding and  
13 succeeding paragraph as though fully set forth herein.

14 237. Plaintiffs and the Class offer code under various Licenses, the most common of which are  
15 set forth in Appendix A. Use of each of the Licensed Materials is allowed only pursuant to the terms of  
16 the applicable Suggested License.

17 238. Plaintiffs and the Class granted Defendants a license to copy, distribute, and/or create  
18 Derivative Works under the Suggested Licenses. Each of the Suggested Licenses requires at least (1) that  
19 attribution be given to the owner of the Licensed Materials used, (2) inclusion of a copyright notice for the  
20 Licensed Materials used, and (3) inclusion of the terms of the applicable Suggested License. When  
21 providing Output, Copilot does not comply with any of these terms.

22 239. Defendants accepted the terms of Plaintiffs' and the Class's Licenses when it used the  
23 licensed code to create Copilot and when it incorporated the licensed code into Copilot. They have  
24 accepted and continue to accept the applicable Licenses every time Copilot Output's Plaintiffs' or the  
25 Class's copyrighted code. As such, contracts have been formed between Defendants on the one hand and  
26 Plaintiffs and the Class on the other.

27 240. Plaintiffs and the Class have performed each of the conditions, covenants, and obligations  
28 imposed on them by the terms of the License associated with their Licensed Materials.

1           241. Plaintiffs and members of the Class hold the copyright in the contents of one or more code  
2 repositories that have been hosted on GitHub’s platform.

3           242. Plaintiffs and the Class have appended one of the Suggested Licenses to each of the  
4 Licensed Materials.

5           243. Plaintiffs and the Class did not know about, authorize, approve, or license the Defendants’  
6 use of the Licensed Materials in the matter at issue in this Complaint before they were used by  
7 Defendants.

8           244. Defendants have substantially and materially breached the applicable Licenses by failing to  
9 provide the source code of Copilot nor a written offer to provide the source code upon the request of each  
10 licensee.

11           245. Defendants have substantially and materially breached the applicable Licenses by failing to  
12 provide attribution to the creator and/or owner of the Licensed Materials.

13           246. Defendants have substantially and materially breached the applicable Licenses by failing to  
14 include copyright notices when Copilot Outputs copyrighted OS code.

15           247. Defendants have substantially and materially breached the applicable Licenses by failing to  
16 identify the License applicable to the Work and/or including its text when Copilot Outputs code including  
17 a portion of a Work.

18           248. Plaintiffs and the Class have suffered monetary damages as a result of Defendants’  
19 conduct.

20           249. The conduct of Defendants is causing and, unless enjoined and restrained by this Court,  
21 will continue to cause Plaintiffs and the Class great and irreparable injury that cannot fully be  
22 compensated or measured in money.

23           250. As a direct and proximate result of these material breaches by Defendants, Plaintiffs and  
24 the Class are entitled to an injunction requiring Defendants to comply with all the terms of any License  
25 governing use of code that was used to train Copilot, otherwise incorporated into Copilot, and/or  
26 reproduced as Output by Copilot.

27           251. Plaintiffs and the Class are further entitled to recover from Defendants the damages  
28 Plaintiffs and the Class sustained—including consequential damages—for Plaintiffs’ and the Class’s costs

1 in enforcing their contractual rights. Plaintiffs and the Class are also entitled to recover as restitution from  
2 Defendants for any unjust enrichment, including gains, profits, and advantages that Defendants have  
3 obtained as a result of their breach of contract.

4  
5 **COUNT 3**  
6 **BREACH OF CONTRACT — SELLING LICENSED MATERIALS**  
7 **IN VIOLATION OF GITHUB’S POLICIES**  
8 **California Common Law**  
9 **(Against GitHub)**

10 252. Plaintiffs and the Class hereby repeat and incorporate by reference each preceding and  
11 succeeding paragraph as though fully set forth herein.

12 253. GitHub’s Privacy Statement, Terms of Service, and GitHub Copilot Terms share  
13 definitions and refer to each other. As such, they are collectively referred to herein as “GitHub’s Policies”  
14 unless a distinction is necessary and are attached as Exhibit 1.

15 254. Plaintiffs and the Class are GitHub users who have accepted GitHub’s Policies. As a result,  
16 Plaintiffs and the Class have formed a contract with GitHub.

17 255. Plaintiffs and the Class have performed each of the conditions, covenants, and obligations  
18 imposed on them by the terms of GitHub’s Policies.

19 256. GitHub’s Policies contain multiple explicit provisions that GitHub will not sell the  
20 Licensed Materials of the Plaintiffs and Class. GitHub’s Terms of Service document provides that the  
21 “License Grant to [GitHub] . . . does not grant GitHub the right to sell Your Content.” Similarly,  
22 GitHub’s Privacy Statement defines “personal data” to include “any . . . documents, or other files”, a  
23 definition that necessarily comprises source code, and hence the Licensed Materials. (As of May 2023,  
24 GitHub has updated this provision on its website to explicitly read “any code, text, . . . documents, or other  
25 files”). Elsewhere, the Privacy Statement provides “We do not sell your personal information,” “No  
26 selling of personal data,” “We *do not* sell your personal data for monetary or other consideration.”  
27 (Emphasis in original).

28 257. By making the Licensed Materials available through Copilot in violation of the Suggested  
Licenses, and charging subscription fees, GitHub has been selling Licensed Materials. By selling the  
Licensed Materials, GitHub has breached these provisions in GitHub’s Policies against selling user data.



- d) An award of damages for harms resulting from Defendants' breach of Licenses;
  - e) An award of damages in the amount Defendants have been unjustly enriched through their conduct as alleged herein as well as punitive damages in connection with this conduct;
  - f) An award of damages for harms resulting from GitHub's breach of the GitHub Policies;
- and

264. Injunctive relief sufficient to alleviate and stop Defendants' unlawful conduct alleged herein.

265. Plaintiffs and the Class are entitled to prejudgment and post-judgment interest on the damages awarded them, 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;

266. 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.

267. Plaintiffs and the Class receive such other or further relief as may be just and proper.

#### **X. JURY TRIAL DEMANDED**

Pursuant to Federal Rule of Civil Procedure 38(b), Plaintiffs demand a trial by jury of all the claims asserted in this Complaint so triable.

1 Dated: January 24, 2024

By:           /s/ Joseph R. Saveri            
Joseph R. Saveri

2 Joseph R. Saveri (State Bar No. 130064)  
3 Cadio Zirpoli (State Bar No. 179108)  
4 Christopher K.L. Young (State Bar No. 318371)  
5 Louis A. Kessler (State Bar No. 243703)  
6 Elissa A. Buchanan (State Bar No. 249996)  
7 Travis Manfredi (State Bar No. 281779)  
8 William W. Castillo Guardado (State Bar No. 294159)  
9 Holden J. Benon (State Bar No. 325847)  
10 **JOSEPH SAVERI LAW FIRM, LLP**  
11 601 California Street, Suite 1000  
12 San Francisco, California 94108  
13 Telephone: (415) 500-6800  
14 Facsimile: (415) 395-9940  
15 Email: jsaveri@saverilawfirm.com  
16 czirpoli@saverilawfirm.com  
17 cyoung@saverilawfirm.com  
18 lkessler@saverilawfirm.com  
19 eabuchanan@saverilawfirm.com  
20 tmanfredi@saverilawfirm.com  
21 wcastillo@saverilawfirm.com  
22 hbenon@saverilawfirm.com

23 Matthew Butterick (State Bar No. 250953)  
24 1920 Hillhurst Avenue, #406  
25 Los Angeles, CA 90027  
26 Telephone: (323) 968-2632  
27 Facsimile: (415) 395-9940  
28 Email: mb@buttericklaw.com

*Counsel for Individual and Representative  
Plaintiffs and the Proposed Class*