

1 Joseph R. Saveri (State Bar No. 130064)  
 2 Cadio Zirpoli (State Bar No. 179108)  
 3 Travis Manfredi (State Bar No. 281779)  
 JOSEPH SAVERI LAW FIRM, LLP  
 4 601 California Street, Suite 1000  
 San Francisco, California 94108  
 Telephone: (415) 500-6800  
 5 Facsimile: (415) 395-9940  
 6 Email: jsaveri@saverilawfirm.com  
 czirpoli@saverilawfirm.com  
 7 tmanfredi@saverilawfirm.com

8 Matthew Butterick (State Bar No. 250953)  
 9 1920 Hillhurst Avenue, #406  
 Los Angeles, CA 90027  
 Telephone: (323) 968-2632  
 Facsimile: (415) 395-9940  
 11 Email: mb@buttericklaw.com

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

14 **UNITED STATES DISTRICT COURT**  
 15 **NORTHERN DISTRICT OF CALIFORNIA**  
 16 **SAN FRANCISCO DIVISION**

17 J. DOE 1 and J. DOE 2, individually and on  
 behalf of all others similarly situated,  
 18 Individual and Representative Plaintiffs,  
 19 v.

Case No.  
**COMPLAINT**  
**CLASS ACTION**

20 GITHUB, INC., a Delaware corporation;  
 21 MICROSOFT CORPORATION, a Washington  
 corporation; OPENAI, INC., a Delaware  
 22 nonprofit corporation; OPENAI, L.P., a  
 Delaware limited partnership; OPENAI GP,  
 23 L.L.C., a Delaware limited liability company;  
 OPENAI STARTUP FUND GP I, L.L.C., a  
 24 Delaware limited liability company; OPENAI  
 STARTUP FUND I, L.P., a Delaware limited  
 25 partnership; OPENAI STARTUP FUND  
 MANAGEMENT, LLC, a Delaware limited  
 26 liability company,  
 27

**DEMAND FOR JURY TRIAL**

Defendants.

**TABLE OF CONTENTS**

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

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

II. JURISDICTION AND VENUE..... 4

III. INTRADISTRICT ASSIGNMENT ..... 4

IV. PARTIES..... 4

    Plaintiffs ..... 4

    Defendants ..... 5

V. AGENTS AND CO-CONSPIRATORS ..... 7

VI. CLASS ALLEGATIONS ..... 8

    A. Class Definitions..... 8

    B. Numerosity..... 9

    C. Typicality..... 9

    D. Commonality & Predominance..... 9

        1. DMCA Violations .....10

        2. Contract-Related Conduct .....10

        3. Unlawful-Competition Conduct .....10

        4. Privacy Violations .....10

        5. Injunctive Relief..... 11

        6. Defenses ..... 11

    E. Adequacy..... 11

    F. Other Class Considerations ..... 11

VII. FACTUAL ALLEGATIONS .....12

    A. Introduction.....12

    B. Codex Outputs Copyrighted Materials Without Following the Terms of  
the Applicable Licenses ..... 13

    C. Copilot Outputs Copyrighted Materials Without Following the Terms of  
the Applicable Licenses .....18

    D. Codex and Copilot Were Trained on Copyrighted Materials Offered Under  
Licenses.....21

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 ..... 22

F. Open-Source Licenses Began to Appear in the Early 1990s ..... 24

G. Microsoft Has a History of Flouting Open-Source License Requirements ..... 26

H. GitHub Was Designed to Cater to Open-Source Projects ..... 28

I. OpenAI Is Intertwined with Microsoft and GitHub..... 30

J. Conclusion of Factual Allegations ..... 32

VIII. CLAIMS FOR RELIEF.....33

IX. DEMAND FOR JUDGMENT ..... 50

X. JURY TRIAL DEMANDED ..... 52

1 Plaintiffs J. Doe 1 and J. Doe 2 (“Plaintiffs”), on behalf of themselves and all others  
2 similarly situated, bring this Class Action Complaint (the “Complaint”) against Defendants  
3 GitHub, Inc.; Microsoft Corporation; OpenAI, Inc.; OpenAI, L.P.; OpenAI GP, L.L.C.; OpenAI  
4 Startup Fund GP I, L.L.C.; OpenAI Startup Fund I, L.P.; and OpenAI Startup Fund  
5 Management, LLC<sup>1</sup> for violation of the Digital Millennium Copyright Act, 17 U.S.C. §§ 1201–  
6 1205 (the “DMCA”); violation of the Lanham Act, 15 U.S.C. § 1125; violation of Unfair  
7 Competition law, *Cal. Bus. & Prof. Code* §§ 17200, *et seq.*; violation of the California Consumer  
8 Privacy Act, *Cal. Civ. Code* § 1798.150 (the “CCPA”); and Breach of Contract regarding the  
9 Suggested Licenses, GitHub’s Privacy Statement, and GitHub’s Terms of Service, *Cal. Bus. &*  
10 *Prof. Code* §§ 22575–22579, *Cal. Civ. Code* § 1798.150. Plaintiffs and the Class also bring this  
11 Complaint against Defendants for their Tortious Interference in Plaintiffs’ Contractual  
12 Relationships; Fraud, and Negligence regarding handling of sensitive data.

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

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

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

---

25 <sup>1</sup> GitHub, Inc. is referred to as “GitHub.” Microsoft Corporation is referred to as “Microsoft.”  
26 OpenAI, Inc.; OpenAI, L.P.; OpenAI GP, L.L.C.; OpenAI Startup Fund GP I, L.L.C.; OpenAI  
27 Startup Fund I, L.P.; and OpenAI Startup Fund Management, LLC are referred to collectively  
28 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.; and  
OpenAI Startup Fund Management, LLC are referred to herein as “Defendants.”

1           3.       GitHub is a company founded in 2008 by a team of open-source enthusiasts. At  
2 the time, GitHub’s stated goal was to support open-source development, especially by hosting  
3 open-source source code on the website github.com. Over the next 10 years, GitHub, based on  
4 these representations succeeded wildly, attracting nearly 25 million developers.

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

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

13           6.       Microsoft invested \$1 billion in OpenAI LP in July 2019 at a \$20 billion valuation.  
14 In 2020, Microsoft became exclusive licensee of OpenAI’s GPT-3 language model—despite  
15 OpenAI’s continued claims its products are meant to benefit “humanity” at large. In 2021,  
16 Microsoft began offering GPT-3 through its Azure cloud-computing platform. On October 20,  
17 2022, it was reported that OpenAI “is in advanced talks to raise more funding from Microsoft” at  
18 that same \$20 billion valuation. Copilot runs on Microsoft’s Azure platform. Microsoft has used  
19 Copilot to promote Azure’s processing power, particularly regarding AI.

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

24           8.       In June 2021, GitHub and OpenAI launched Copilot, an AI-based product that  
25 promises to assist software coders by providing or filling in blocks of code using AI. GitHub  
26 charges Copilot users \$10 per month or \$100 per year for this service. Copilot ignores, violates,  
27 and removes the Licenses offered by thousands—possibly millions—of software developers,  
28 thereby accomplishing software piracy on an unprecedented scale. Copilot outputs text derived

1 from Plaintiffs' and the Class's Licensed Materials without adhering to the applicable License  
2 Terms and applicable laws. Copilot's output is referred herein as "Output."

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

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

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

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

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

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

21 15. It is not fair, permitted, or justified. On the contrary, Copilot's goal is to replace a  
22 huge swath of open source by taking it and keeping it inside a GitHub-controlled paywall. It  
23 violates the licenses that open-source programmers chose and monetizes their code despite  
24 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  
26 that it interprets using given criteria. Feedback is then given to it to fine-tune its Output until it  
27 can provide 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.

## II. JURISDICTION AND VENUE

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

17. Jurisdiction and venue are proper in this judicial district pursuant to Defendants' violation of the Digital Millennium Copyright Act, 17 U.S.C. §§ 1201-1205 (the "DMCA"); Reverse Passing Off, Unjust Enrichment, and Unfair Competition under the Lanham Act, 15 U.S.C. § 1125; and because a substantial part of the events giving rise to Plaintiff's claims occurred in this District, Plaintiff J. Doe 1 resides in California, a substantial portion of the affected interstate trade and commerce was carried out in this District, and three or more of the Defendants reside in this District and/or are licensed to do business in this District. Each Defendant has transacted business, maintained substantial contacts, and/or committed overt acts in furtherance of the illegal scheme and conspiracy throughout the United States, including in this District. Defendants' conduct has had the intended and foreseeable effect of causing injury to persons residing in, located in, or doing business throughout the United States, including in this District.

## III. INTRADISTRICT ASSIGNMENT

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

## IV. PARTIES

### PLAINTIFFS

19. Plaintiff J. Doe 1 is a resident of the State of California. Plaintiff Doe 1 published Licensed Materials they owned a copyright interest in to at least one GitHub repository under one of the Suggested Licenses. Specifically, Doe 1 has published Licensed Materials they claim a

1 copyright interest in under the following Suggested Licenses: MIT License and GNU General  
2 Public License version 3.0. Plaintiff was, and continues to be, injured during the Class Period as a  
3 result of Defendants’ unlawful conduct alleged herein.

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

#### 11 **DEFENDANTS**

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

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

25 23. Defendant OpenAI, Inc. is a Delaware nonprofit corporation with its principal  
26 place of business located at 3180 18th Street, San Francisco, CA 94110. OpenAI, Inc. is a party to  
27 the unlawful conduct alleged herein. It—along with OpenAI, L.P.—programed, trained, and  
28 maintains Codex, which infringes all the same rights at Copilot and is also an integral piece of



1 Copilot. Copilot requires Codex to function. OpenAI, Inc. is a party to the unlawful conduct  
2 alleged herein. OpenAI, Inc. founded, owns, and exercises control over all the other OpenAI  
3 entities, including those set forth in Paragraphs 24–28.

4 24. Defendant OpenAI, L.P. is a Delaware limited partnership with its principal place  
5 of business located at 3180 18th Street, San Francisco, CA 94110. OpenAI, L.P. is a party to the  
6 unlawful conduct alleged herein. Its primary activity is research and technology. OpenAI, L.P. is a  
7 wholly owned subsidiary of OpenAI, Inc. that is operated for profit. OpenAI, L.P. is the OpenAI  
8 entity that co-created Copilot and offers it jointly with GitHub. OpenAI’s revenue, including  
9 revenue from Copilot, is received by OpenAI, L.P. OpenAI, Inc. controls OpenAI, L.P. directly  
10 and through the other OpenAI entities.

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

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





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

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

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

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

16 **B. Numerosity**

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

21 **C. Typicality**

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

27 **D. Commonality & Predominance**

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

1           **1.     DMCA Violations**

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

6           **2.     Contract-Related Conduct**

- 7           • Whether Defendants violated the Licenses governing use of the Licensed  
8           Materials by using them to train Copilot and for republishing those materials  
9           without appending the required Attribution, Copyright Notice, or License  
10          Terms.  
11          • Whether Defendants interfered in contractual relations between the Class and  
12          the public regarding the Licensed Materials by concealing the License Terms.  
13          • Whether GitHub committed Fraud when it promised not to sell or distribute  
14          Licensed Materials outside GitHub in the GitHub Terms of Service and  
15          Privacy Statement.

16          **3.     Unlawful-Competition Conduct**

- 17          • Whether Defendants passed-off the Licensed Materials as its own creation  
18          and/or Copilot’s creation.  
19          • Whether Defendants were unjustly enriched by the unlawful conduct alleged  
20          herein.  
21          • Whether Defendants Copilot-related conduct constitutes Unfair Competition  
22          under California law.

23          **4.     Privacy Violations**

- 24          • Whether GitHub violated the Class’s rights under the California Consumer  
25          Privacy Act (“CCPA”), the GitHub Privacy Statement, and/or the California  
26          Constitution by, *inter alia*, sharing the Class’s sensitive personal information  
27          (or, in the alternative, by not addressing an ongoing data breach of which it is  
28          aware); creating a product that contains personal data GitHub cannot delete,

1 alter, nor share with the applicable Class member; and selling the Class's  
2 personal data.

- 3 • Whether GitHub committed Negligence when it failed to stop a still-ongoing  
4 data breach it was and continues to be aware of.

5 **5. Injunctive Relief**

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

8 **6. Defenses**

- 9 • Whether any affirmative defense excuses Defendants' conduct.  
10 • Whether any statutes of limitation limit Plaintiffs' and the Class's potential for  
11 recovery.  
12 • Whether any applicable statutes of limitation should be tolled as a result of  
13 Defendants' fraudulent concealment of their unlawful conduct.

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

16 **E. Adequacy**

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

23 **F. Other Class Considerations**

24 40. Defendants have acted on grounds generally applicable to the Class, thereby  
25 making final injunctive relief appropriate with respect to the Class as a whole.

26 41. This class action is superior to alternatives, if any, for the fair and efficient  
27 adjudication of this controversy. Prosecuting the claims pleaded herein as a class action will  
28

1 eliminate the possibility of repetitive litigation. There will be no material difficulty in the  
2 management of this action as a class action.

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

## 6 VII. FACTUAL ALLEGATIONS

### 7 A. Introduction

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

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

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

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

22 47. Codex and Copilot are related. Copilot is a joint venture between GitHub and  
23 OpenAI. Codex is a standalone product released by OpenAI that also “powers GitHub Copilot,  
24 which [OpenAI] built and launched in partnership with GitHub.”<sup>5</sup> “GitHub Copilot uses the  
25 OpenAI Codex to suggest code and entire functions in real-time, right from your editor.”<sup>6</sup> The  
26

---

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

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

1 version of Codex used to power Copilot is modified from OpenAI’s standalone version.  
2 Defendants have kept secret the details of Codex’s modifications and its integration into or  
3 interaction with Copilot.

4 **B. Codex Outputs Copyrighted Materials Without Following the Terms of the**  
5 **Applicable Licenses**

6 48. Below is an explanation of how Codex functions. When Codex is prompted with:

7 `function isEven(n) {`

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

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28



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

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

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

---

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

1           51.     The penultimate line indicates “`isEven(-1)`” should return “??”. This is an  
2 error, as “`isEven(-1)`” should return “false”.

3           52.     Codex cannot and does not understand the meaning of software code or any other  
4 Licensed Materials. But in training, what became Codex was exposed to an enormous amount of  
5 existing software code (its “Training Data”) and—with input from its trainers and its own  
6 internal processes—inferred certain statistical patterns governing the structure of code and other  
7 Licensed Materials. The finished version of Codex, once trained, is known as a “Model.”

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

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

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

26           56.     Codex does not identify the owner of the copyright to this Output, nor any  
27 other—it has not been trained to provide Attribution. Nor does it include a Copyright Notice nor  
28 any License Terms attached to the Output. This is by design—Codex was not coded or trained to

1 track or reproduce such data. The Output in the example above is taken from *Eloquent Javascript*  
2 by Marijn Haverbeke.<sup>8</sup>

3 57. Here is the exercise from *Eloquent Javascript*:

```
4 // Your code here.  
5  
6 console.log(isEven(50));  
7 // → true  
8 console.log(isEven(75));  
9 // → false  
10 console.log(isEven(-1));  
11 // → ??
```

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

17 59. Haverbeke provides the following solution to the function discussed above:

```
18 function isEven(n) {  
19     if (n == 0) return true;  
20     else if (n == 1) return false;  
21     else if (n < 0) return isEven(-n);  
22     else return isEven(n - 2);  
23 }  
24  
25 console.log(isEven(50));
```

---

26 <sup>8</sup> <https://eloquentjavascript.net/code/#3.2>. *Eloquent Javascript* is “Licensed under a Creative  
27 Commons [A]tribution-[N]oncommercial license. All code in this book may also be considered  
28 licensed under an 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 // → true
2 console.log(isEven(75));
3 // → false
4 console.log(isEven(-1));
5 // → false
```

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

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

13 62. The MIT license provides that “The above copyright notice and this permission  
14 notice shall be included in all copies or substantial portions of the Software.”<sup>9</sup> Any person taking  
15 this code directly from *Eloquent JavaScript* would have direct access to these License Terms and  
16 know to follow them if incorporating the Licensed Materials into a derivative work and/or  
17 copying them. Codex does not provide these License Terms.

18 63. OpenAI Codex’s Output would frequently, perhaps even constantly, contain  
19 Licensed Materials, i.e., it would have conditions associated with it through its associated license.  
20 In its 2021 research paper about Codex called “Evaluating Large Language Models Trained on  
21 Code,” OpenAI stated Codex’s Output is “often incorrect” and can contain security  
22 vulnerabilities and other “misalignments” (meaning, departures from what the user requested).

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

---

27  
28 <sup>9</sup> See Appendix A for full text of the MIT License.

1           65.     Ultimately, Codex derives its value primarily from its ability to locate and output  
2 potentially useful Licensed Materials. And from its obfuscation of any rights associated with  
3 those materials.

4     **C.     Copilot Outputs Copyrighted Materials Without Following the Terms of the**  
5     **Applicable Licenses**

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

8           67.     Copilot is installed by the end user as an extension to various code editors,  
9 including Microsoft’s Visual Studio and VS Code. As the user types into the editor, their code is  
10 uploaded in real time to Microsoft’s Azure cloud platform, where they become prompts for  
11 Copilot.

12           68.     When we give Copilot the same prompt discussed above in Paragraph 48,  
13 “function isEven(n) {”, it interprets the prompt as the beginning of a function written in  
14 the JavaScript language that will test whether a number is even, same as Codex.

15           69.     However, the Output of Copilot in response to the prompt is different than Codex,  
16 namely:

```
17     function isEven(n) {  
18         return n % 2 === 0;  
19     }
```

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

23           71.     Copilot’s Output, like Codex’s, is derived from existing code. Namely, sample  
24 code that appears in the online book *Mastering JS*, written by Valeri Karpov.<sup>10</sup> Like Codex’s  
25 Output, Copilot’s is also based upon copyrighted educational material. *Mastering JS* is a set of  
26 educational exercises for programmers. Like *Eloquent Javascript*, there are many copies of  
27

---

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

1 Karpov’s exercise stored in public repositories on GitHub. Programmers working through  
2 *Mastering JS* store their answers there.

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

```
6     function isPrime(n) {  
7         if (n < 2) {  
8             return false;  
9         }  
10        for (let i = 2; i < n; i++) {  
11            if (n % i === 0) {  
12                return false;  
13            }  
14        }  
15        return true;  
16    }
```

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

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

---

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

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

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

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

22 77. As with Codex, Copilot does not provide the end user any attribution of the  
23 original author of the code, nor anything about their license requirements. There is no way for the  
24 Copilot user to know that they must provide attribution, copyright notice, nor a copy of the  
25 license’s text. And with regard to the GNU Free Documentation License, Copilot users would  
26 not be aware that they are limited in what conditions they can place on the use of derivative works  
27

---

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

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

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

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

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

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

22 81. The words “study” and “training” and “learning” in connection with AI describe  
23 algorithmic processes that are not analogous to human reasoning. An AI models cannot “learn”  
24 as humans do, nor can it “understand” semantics and context the way humans do. Rather, it  
25 detects statistically significant patterns in its training data and provides Output derived from its  
26 training data when statistically appropriate. A “brute force” approach like this would not be  
27 efficient nor even possible for humans. A human could not memorize, statistically analyze, and  
28 easily access thousands of gigabytes of existing code, a task now possible for powerful computers



1 like those that make up Microsoft’s Azure cloud platform. To accomplish the same task, a human  
2 may search for Licensed Materials that serve their purpose if they believe such materials exist.  
3 And if that human finds such materials, they will probably abide by its License Terms rather than  
4 risk infringing its owners’ rights. At the very least, if they incorporate those Licensed Materials  
5 into their own project without following its terms they will be doing so knowingly.

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

7 82. GitHub and OpenAI have not provided much detail regarding what data Codex  
8 and OpenAI were trained on. Plaintiffs know for certain from GitHub and OpenAI’s statements,  
9 that both systems were trained on publicly available GitHub repositories, with Copilot having  
10 been trained on all available public GitHub repositories. Thus, if Licensed Materials have been  
11 posted to a GitHub public repository, Plaintiffs and the Class can be reasonably certain it was  
12 ingested by Copilot and is sometimes returned to users as Output.

13 83. According to OpenAI, Codex was trained on “billions of lines of source code from  
14 publicly available sources, including code in public GitHub repositories”. Similarly, GitHub has  
15 described<sup>13</sup> Copilot’s training material as “billions of lines of public code.” GitHub researcher  
16 Eddie Aftandilian confirmed in a recent podcast<sup>14</sup> that Copilot is “train[ed] on public repos on  
17 GitHub.”

18 84. In a recent customer-support message, GitHub’s support department clarified  
19 certain facts about training Copilot. First, GitHub said that “training for Codex (the model used  
20 by Copilot) is done by OpenAI, not GitHub.” Second, in its support message, GitHub put  
21 forward a more detailed justification for its use of copyrighted code as training data:

---

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

12 The claim that training ML models on publicly available code is widely accepted as fair use is not  
13 true. And regardless of this concept’s level of acceptance in “the machine learning community,”  
14 under Federal law, it is illegal.

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

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

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

---

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

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

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

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

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

6 90. GitHub concedes that in ordinary use, Copilot will reproduce passages of code  
7 verbatim: “Our latest internal research shows that about 1% of the time, a suggestion [Output]  
8 may contain some code snippets longer than ~150 characters that matches” code from the  
9 training data. This standard is more limited than is necessary for copyright infringement. But  
10 even using GitHub’s own metric and the most conservative possible criteria, Copilot has violated  
11 the DMCA at least tens of thousands of times.

12 91. In June 2022, Copilot had 1,200,000 users. If only 1% of users have ever received  
13 Output based on Licensed Materials and only once each, Defendants have “only” breached  
14 Plaintiffs’ and the Class’s Licenses 12,000 times. However, each time Copilot outputs Licensed  
15 Materials without attribution, the copyright notice, or the License Terms it violates the DMCA  
16 three times. Thus, even using this extreme underestimate, Copilot has “only” violated the  
17 DMCA 36,000 times.<sup>20</sup> Because Copilot constantly Outputs code as a user writes, and because  
18 nearly all of Copilot’s training data was Licensed Material, this number is most likely  
19 exponentially lower than the true number of breaches and DMCA violations.

## 20 **F. Open-Source Licenses Began to Appear in the Early 1990s**

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

---

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

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

26 <sup>20</sup> These violations of Section 1202 of the DMCA each incur statutory damages of “not less than  
27 \$2,500 or more than \$25,000.” 17 U.S.C. § 1203(c)(3)(B). This extremely conservative estimate  
28 of Defendants’ number of direct violations translates to \$90 million to \$900 million in statutory  
damages.

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

4           94. The GPL is a software license. But whereas most software licenses required  
5 payment, software under the GPL is provided for free. Whereas most software licenses did not  
6 include source code, GPL software always included source code. And whereas most software  
7 licenses prohibited derivative works, the GPL not only allowed it, but encouraged it.

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

12           96. One license requirement was that a program derived from GPL software had to  
13 redistribute certain information about that software:

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

19 Failure to adhere to these conditions constituted a violation of the license, triggering the  
20 possibility of legal action. Provisions of the GPL are enforceable, and many GPL licensors have  
21 sought to enforce GPL licenses through court proceedings and other litigation.

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

---

27  
28 <sup>21</sup> <https://www.gnu.org/licenses/old-licenses/gpl-1.0.en.html>.

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

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

6 **G.     Microsoft Has a History of Flouting Open-Source License Requirements**

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

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

15           102.    Linux represented a competitive threat to Windows. It ran on the same hardware.  
16 It performed many of the same functions. It was free. Many programmers at the time considered  
17 Linux to be functionally superior to Windows.

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

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

24           105.    In 1998, a source at Microsoft leaked what became known as the “Halloween  
25 Documents”, revealing Microsoft’s thinking on how to counter the competitive threat from  
26 Linux. Among other things, the documents emphasized the importance of countering the “long  
27  
28

1 term developer mindshare threat”, and concluded that to defeat open source, “[Microsoft] must  
2 target a process rather than a company”.<sup>22</sup>

3 106. In 2001, Microsoft CEO Steve Ballmer said “The way the [GPL] is written, if you  
4 use any open-source software, you must make the rest of your software open source. . . . Linux is  
5 a cancer that attaches itself in an intellectual property sense to everything it touches.”<sup>23</sup>

6 Ballmer’s summary of GPL licensing was not accurate. In 2001, Linux was being used by  
7 corporations of every size. The growth of open source up to that point, and since, has been made  
8 possible by the open-source community’s respect for and compliance with applicable licenses.

9 107. In 2001, Microsoft was the defendant in a major software-related antitrust case,  
10 *United States v. Microsoft Corporation*.<sup>24</sup> In this case, the U.S. Department of Justice accused  
11 Microsoft of maintaining a software monopoly by illegally imposing technical restrictions on  
12 manufacturers of personal computers, including “tying” violations related to the Internet  
13 Explorer web browser. Judge Thomas Penfield Jackson, who presided over the antitrust trial,  
14 opined that Microsoft is “a company with an institutional disdain for both the truth and for rules  
15 of law that lesser entities must respect. It is also a company whose ‘senior management’ is not  
16 averse to offering specious testimony to support spurious defenses to claims of its wrongdoing.”<sup>25</sup>

17 108. In 2007, Microsoft admitted that it tried to influence the vote of an ISO open-  
18 standards committee by offering money to certain business partners in Sweden to vote for  
19 Microsoft’s preferred outcome.<sup>26</sup>

20 109. After observing the rapid growth of Amazon’s original cloud computing products,  
21 Microsoft has expanded its business into cloud computing, which it has branded Microsoft Azure  
22 or simply Azure. Microsoft announced Azure to developers in 2008. It was formally released in  
23

---

24 <sup>22</sup> <http://www.catb.org/esr/halloween/halloween1.html>.

25 <sup>23</sup> <https://lwn.net/2001/0607/a/esr-big-lie.php3>.

26 <sup>24</sup> No. Civ.A. 00-1457 TPJ.

27 <sup>25</sup> *Jackson v. Microsoft Corp.*, 135 F. Supp. 2d 38 (D.D.C. 2001).

28 <sup>26</sup> <https://learn.microsoft.com/en-us/archive/blogs/jasonmatusow/open-xml-the-vote-in-sweden/>.

1 2010. Azure uses large-scale virtualization at Microsoft data centers and offers many hundreds of  
2 services, including infrastructure as a service (“IaaS”), platform as a service (“PaaS”), compute  
3 services, Azure Active Directory, mobile services, storage services, communication services, data  
4 management, messaging, developer services, Azure AI, blockchain, and others.

#### 5 **H. GitHub Was Designed to Cater to Open-Source Projects**

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

9 111. Linus Torvalds set about writing a new source-control system. He named his new  
10 system Git. He released it under the GPL. It quickly became the source-control system of choice  
11 for open-source programmers.

12 112. A single software project stored in Git is called a *source repository*, commonly  
13 shortened to *repository* or just *repo*. A Git source repository would typically be stored on a  
14 networked server accessible to a group of programmers.

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

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

24 115. Most open-source programmers used GitHub to create “public” repositories,  
25 meaning that anyone could view them & access them. GitHub also allowed programmers and  
26 organizations to create “private” repositories, which were not accessible from the public GitHub  
27 website, and required password access.

28

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

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

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

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

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

23           121. In October 2018, Microsoft acquired GitHub for \$7.5 billion. It was important to  
24 Microsoft that programmers use GitHub. Microsoft had developed a well-deserved poor  
25 reputation because of its documented vaporware, FUD, and other business practices, including  
26 those targeted at open-source programs and programming, and open-source licensing specifically.  
27 Microsoft made false and misleading statements and omissions to assuage such concerns,  
28



1 including its primary mantra intended to win over the open-source community: “Microsoft Loves  
2 Open Source.”

3 **I. OpenAI Is Intertwined with Microsoft and GitHub**

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

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

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

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

19 126. OpenAI, Inc. obtained its initial source of training data from its founders’  
20 companies. According to reporting at the time, Musk and Altman planned to “pool[] online data  
21 from their respective companies” to serve as training data for OpenAI, Inc. projects. Musk  
22 planned to contribute data from Tesla; Altman planned to have Y Combinator companies “share  
23 their data with OpenAI.”<sup>27</sup>

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

---

27 <sup>27</sup> [https://www.wired.com/2015/12/elon-musks-billion-dollar-ai-plan-is-about-far-more-than-](https://www.wired.com/2015/12/elon-musks-billion-dollar-ai-plan-is-about-far-more-than-saving-the-world/)  
28 [saving-the-world/](https://www.wired.com/2015/12/elon-musks-billion-dollar-ai-plan-is-about-far-more-than-saving-the-world/).

1           128. In July 2019, OpenAI, L.P. accepted a \$1 billion investment from Microsoft. In  
2 addition to cash, Microsoft would become the exclusive licensor of certain OpenAI, LP products  
3 (including GPT-3, described below in Paragraph 131). Also, as part of this alliance, OpenAI, LP  
4 would use Microsoft’s cloud-computing platform, Azure, exclusively to develop and host its  
5 products. Some portion of Microsoft’s investment was paid in credits for use of Azure rather  
6 than cash. Finally, Microsoft and OpenAI agreed to “jointly build new Azure AI supercomputing  
7 technologies.”

8           129. Azure is a major growth area for Microsoft. In its most recent earnings report on  
9 October 25, 2022, “Azure and other cloud services” grew by 35% from the previous quarter, more  
10 than any other product.<sup>28</sup> Azure has grown rapidly since Microsoft began its partnership with  
11 OpenAI in 2016. Its revenue grew by 50% or more every quarter from 2016 through the first three  
12 quarters of 2020.

13           130. In May 2020, Microsoft and OpenAI announced they had jointly built a  
14 supercomputer in Azure that would be used exclusively by OpenAI to train its AI models.  
15 Microsoft’s influence over and frequent collaboration with OpenAI has led some to describe  
16 Microsoft as “the unofficial owner of OpenAI.”<sup>29</sup>

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

21           132. Sometime after July 2019, OpenAI and Microsoft began collaborating on a code-  
22 completion product for GitHub that would use Codex as its underlying model. This product  
23 became known as Copilot.

24           133. On September 28, 2022, OpenAI released an image-generation AI called DALL-  
25 E-2. Much like Copilot, DALL-E-2 removes any attribution and/or copyright notice from the  
26

---

27 <sup>28</sup> <https://www.microsoft.com/en-us/Investor/earnings/FY-2023-Q1/press-release-webcast/>.

28 <sup>29</sup> <https://venturebeat.com/ai/what-to-expect-from-openais-codex-api/>.

1 images it uses to create derivative works. Like with Codex, here, OpenAI ignores the rights of the  
2 owners of copyrights to images it has ingested.

3 134. In another joint project, Microsoft and OpenAI recently launched a preview of a  
4 product called “Azure OpenAI Service.”<sup>30</sup> This service will “Leverage large-scale, generative AI  
5 models with deep understandings of language and code to enable new reasoning and  
6 comprehension capabilities for building cutting-edge applications. Apply these coding and  
7 language models to a variety of use cases, such as writing assistance, code generation, and  
8 reasoning over data. Detect and mitigate harmful use with built-in responsible AI and access  
9 enterprise-grade Azure security.”

#### 10 **J. Conclusion of Factual Allegations**

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

15 136. Defendants have made no attempt to comply with the open-source licenses that  
16 are attached to much of their training data. Instead, they have pretended those licenses do not  
17 exist, and trained Codex and Copilot to do the same. By simultaneously violating the open-source  
18 licenses of tens-of-thousands—possibly millions—of software developers, Defendants have  
19 accomplished software piracy on an unprecedented scale. As Microsoft’s Co-Founder Bill Gates  
20 once said regarding software piracy: “the thing you do is theft.”<sup>31</sup>

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

---

27 <sup>30</sup> <https://azure.microsoft.com/en-us/products/cognitive-services/openai-service/>.

28 <sup>31</sup> [https://www.digibarn.com/collections/newsletters/homebrew/V2\\_01/gatesletter.html](https://www.digibarn.com/collections/newsletters/homebrew/V2_01/gatesletter.html)

**VIII. CLAIMS FOR RELIEF**

**COUNT I  
VIOLATION OF THE DIGITAL MILLENIUM COPYRIGHT ACT  
17 U.S.C. §§ 1201–1205  
(Direct, Vicarious, and Contributory)  
(Against All Defendants)**

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

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

140. Plaintiffs and members of the Class included the following Copyright Management Information (as defined in Section 1202(c) of the DMCA) (“CMI”) in the Licensed Materials:

- a. copyright notices;
- b. the title and other information identifying the Licensed Materials;
- c. the name of, and other identifying information about, the authors of the Licensed Materials;
- d. the name of, and other identifying information about, the copyright owners of the Licensed 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.

141. 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.

142. 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.

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

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

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

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

16           147. Defendants had access to but were not licensed by Plaintiffs nor the Class to create  
17 Derivative Works<sup>32</sup> based upon the Licensed Materials.

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

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

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

---

26  
27 <sup>32</sup> "Derivative Works" as used herein refers to Copilot's Output to the extent they are derived  
28 from Licensed Materials. The definition also includes the Copilot product itself, which is a  
Derivative Work based upon a large corpus of Licensed Materials.

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

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

16           153. Defendants removed or altered CMI from Licensed Materials owned by Plaintiffs  
17 and the Class while possessing reasonable grounds to know that it would induce, enable, facilitate,  
18 and/or conceal infringement of copyright in violation of the DMCA. By omitting and concealing  
19 CMI from Copilot's Output, Defendants have reasonable grounds to know that innocent  
20 infringers are induced or enabled to copy the Licensed Materials, because CMI has been  
21 removed. Without the CMI, Defendants have reasonable grounds to know copyright infringement  
22 is facilitated or concealed, because Plaintiffs and the Class have the difficult or impossible task of  
23 proving the Licensed Materials belong to them.

24           154. Defendants knowingly provided CMI that is false with respect to the Licensed  
25 Materials. Defendants have a business practice of asserting and/or implying that Copilot is the  
26 author of the Licensed Materials. Defendants knowingly distributed CMI that is false, with  
27 respect to the Licensed Materials. Defendants have a business practice of asserting and/or  
28 implying that Copilot is the author of the Licensed Materials.

1           155. Defendants provided or distributed false CMI from the Licensed Materials with  
2 respect to Copilot's Output with the intent and foreseeable result to induce, enable, facilitate, or  
3 conceal infringement. Defendants have a business practice of asserting and/or implying that  
4 Copilot is the author of the Licensed Materials. This false CMI induces or enables Defendants or  
5 Copilot users to copy the Licensed Materials. Defendants' false description of the source of  
6 Copilot's Output facilitated or concealed infringement by Defendants and Copilot users because  
7 Plaintiffs and the Class have the difficult or impossible task of proving that the copyrights to the  
8 suggested portions of their Licensed Materials belong to them once those Licensed Materials  
9 have been delinked from all identifying information and all license terms governing their use.

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

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

19           158. Defendants did not request or obtain permission from Plaintiffs and the Class to  
20 use the Licensed Materials for Defendants' Copilot product.

21           159. Defendants use of the Licensed Materials does not follow the requirements of the  
22 Suggested Licenses associated with the Licensed Materials. In particular, Copilot fails to provide  
23 attribution for the creator nor the owner of the Work. Copilot fails to include the required  
24 copyright notice included in the License. Copilot fails to include the applicable Suggested  
25 License's text.

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

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

24           Microsoft Corporation Annual Report, Form 10-K at 27 (July 28, 2022) (emphasis added).<sup>33</sup>

25           161.     GitHub, which offers the Copilot product jointly with OpenAI, also has extensive  
 26           experience with the DMCA. GitHub knows or reasonably should know that the Licensed  
 27           Materials it hosts are subject to copyright. It provides the language of the Suggested Licenses to  
 28           users, all of which include copyright notices. Its 2022 Transparency Report—January to June<sup>34</sup>  
 states: “Copyright-related takedowns (which we often refer to as DMCA takedowns) are  
 particularly relevant to GitHub because so much of our users’ content is software code and can be  
 eligible for copyright protection.”<sup>35</sup> In the first six months of 2022, GitHub processed 1220  
 DMCA takedown requests. Its DMCA Takedown Policy<sup>36</sup> notes “GitHub probably never would  
 have existed without the DMCA.”

---

<sup>33</sup> <https://microsoft.gcs-web.com/static-files/07cf3c30-cfc3-4567-b20f-f4b0f0bd5087/>.

<sup>34</sup> <https://github.blog/2022-08-16-2022-transparency-report-january-to-june/>.

<sup>35</sup> <https://github.blog/2022-08-16-2022-transparency-report-january-to-june/>.

<sup>36</sup> <https://docs.github.com/en/site-policy/content-removal-policies/dmca-takedown-policy#what-is-the-dmca/>.



1           162. GitHub also knows or reasonably should know the portions of the DMCA giving  
2 rise to Plaintiffs’ claim. In its 2021 Transparency Report, “Before removing content based on  
3 alleged circumvention of copyright controls (under Section 1201 of the US DMCA or similar laws  
4 in other countries), we carefully review both the legal and technical claims, and we sponsor a  
5 Developer Defense Fund to provide developers with meaningful access to legal resources.”<sup>37</sup>

6           163. GitHub is aware that Copilot’s removal of CMI is illegal. For example, it states  
7 that “publishing or sharing tools that enable circumvention are not [permitted]”<sup>38</sup> and  
8 “Distributing tools that enable circumvention is prohibited, even if their use by developers falls  
9 under the exemption [for security research].”<sup>39</sup> GitHub has also frequently published articles  
10 discussing the DMCA, its application, and the Copyright Office’s guidance on its scope and  
11 exceptions.<sup>40</sup>

12           164. Unless Defendants are enjoined from violating the DMCA, Plaintiffs and the Class  
13 will suffer great and irreparable harm by depriving them of the right to identify and control the  
14 reproduction and/or distribution of their copyrighted works, to have the terms of their open-  
15 source licenses followed, and to pursue copyright-infringement remedies. Defendants will not be  
16 damaged if they are required to comply with the DMCA. Plaintiffs and the Class members are  
17 therefore entitled to an injunction barring Defendants from violating the DMCA and impounding  
18 any device or product that is in the custody or control of Defendants and that the court has  
19 reasonable cause to believe was involved in a violation of the DMCA.

20           165. Plaintiffs and the Class are further entitled to recover from Defendants the actual  
21 or statutory damages Plaintiffs and the Class sustained pursuant to 17 U.S.C. § 1203(c) and for  
22 Plaintiffs’ and the Class’s costs and attorneys’ fees in enforcing the Licenses. Plaintiffs and the  
23 Class are also entitled to recover as restitution from Defendants for any unjust enrichment,  
24

---

25 <sup>37</sup> <https://github.blog/2022-01-27-2021-transparency-report/>.

26 <sup>38</sup> <https://github.blog/2020-11-19-take-action-dmca-anti-circumvention-and-developer-innovation/#what-dmca-exemptions-do-not-do/>.

27 <sup>39</sup> <https://github.blog/2021-11-23-copyright-office-expands-security-research-rights/>.

28 <sup>40</sup> *See, e.g.*, Footnotes 34–39.

1 including gains, profits, and advantages that Defendants have obtained as a result of their breach  
2 of the Licenses.

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

5 167. Defendants induced Copilot users to unknowingly violate the DMCA by  
6 withholding attribution, licensing, and other information as described herein.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

4 **COUNT III**  
5 **TORTIOUS INTERFERENCE IN A CONTRACTUAL RELATIONSHIP**  
6 **Common Law**  
7 **(Against All Defendants)**

8 184. Plaintiffs and the Class hereby repeat and incorporate by reference each preceding  
9 and succeeding paragraph as though fully set forth herein.

10 185. Defendants have wrongfully interfered with the business interests and  
11 expectations of Plaintiffs and the Class by improperly using Copilot to create Derivative Works  
12 that compete against OSC.

13 186. At GitHub’s upcoming yearly conference, GitHub Universe 2022, it will host a  
14 presentation called “How to compete with open source—and win.”

15 187. Plaintiffs and the Class have suffered monetary, reputational, and other damages  
16 as a result of Defendants’ conduct.

17 188. The harm was the actual, proximate, intentional, direct, and foreseeable  
18 consequence of Defendant’s conduct.

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

22 **COUNT IV**  
23 **FRAUD**  
24 **Common Law**  
25 **(Against GitHub)**

26 190. Plaintiffs and the Class hereby repeat and incorporate by reference each preceding  
27 and succeeding paragraph as though fully set forth herein.

28 191. GitHub made certain representations to Plaintiffs and the Class to induce them to  
publicly post their code on GitHub. Specifically, in both its Terms of Service and its Privacy  
Statement, GitHub promises not to sell Licensed Materials or anything else uploaded to or shared

1 with GitHub. It also promises not to distribute Licensed Materials outside GitHub. As explained  
2 above, Copilot operates on an individual’s computer as an extension to their editor as well as on  
3 Microsoft’s Azure cloud platform. Neither are part of GitHub. It Outputs in the user’s editor,  
4 which is not part of GitHub.

5 192. Plaintiffs and the Class relied upon those representations in choosing to upload  
6 Licensed Materials to GitHub. GitHub has long held itself out as the best place to host open-  
7 source code repositories. It has courted the business of users it expects will include Licenses with  
8 their code. It facilitates this by allowing users to easily select the name of a license, including the  
9 Suggested Licenses, when creating a repository rather than finding the text of the license and  
10 adding it themselves. GitHub provides the terms, it can hardly claim to be unaware of what they  
11 are or what they mean. If it didn’t understand the requirements of a given Suggested License, it  
12 would not have provided it as an option to its users.

13 193. GitHub failed to honor its representations in creating and operating Copilot. It  
14 sells Plaintiffs’ and the Class’s Licensed Materials as part of Copilot. It also distributes them. It  
15 does so without following any of the License Terms.

16 194. As such, GitHub failed to honor its representations in operating Copilot.

17 195. The conduct of GitHub is causing and, unless enjoined and restrained by this  
18 Court, will continue to cause Plaintiffs and the Class great and irreparable injury that cannot fully  
19 be compensated or measured in money. Namely, it will continue the proliferation of copies of  
20 Licensed Materials divorced from their licenses and identifying information until infringement is  
21 so prevalent no amount of enforcement by Plaintiffs and the Class could stop its spread.

22 **COUNT V**  
23 **FALSE DESIGNATION OF ORIGIN—REVERSE PASSING OFF**  
24 **15 U.S.C. § 1125**  
**(GitHub and OpenAI)**

25 196. Plaintiffs and the Class hereby repeat and incorporate by reference each preceding  
26 and succeeding paragraph as though fully set forth herein.

27 197. GitHub and OpenAI have used or made, and will continue to use or make, in  
28 commerce throughout the United States, including in California, one or more words, terms,

1 names, symbols, or devices, or any combination thereof, or any false and/or misleading  
2 designation of origin, false and/or misleading description of fact, or false and/or misleading  
3 representation of fact that is likely to cause consumer confusion, or to cause mistake, or to deceive  
4 as to the affiliation, connection, or association of Plaintiffs' and the Class's Licensed Materials  
5 and Copilot, or as to the origin, sponsorship, or approval of Plaintiffs' and the Class's Licensed  
6 Materials and Copilot.

7 198. As a result, GitHub and OpenAI have intentionally violated 15 U.S.C. §  
8 1125(a)(1)(A).

9 199. As an actual and proximate result of GitHub's and OpenAI's acts, Plaintiffs and  
10 the Class have suffered and continue to suffer harm.

11 **COUNT VI**  
12 **UNJUST ENRICHMENT**  
13 ***Cal. Bus. & Prof. Code §§ 17200, et seq. and Common Law***  
**(GitHub and OpenAI)**

14 200. Plaintiffs and the Class hereby repeat and incorporate by reference each preceding  
15 and succeeding paragraph as though fully set forth herein.

16 201. Plaintiffs and the Class have invested substantial time and energy in creating the  
17 Licensed Materials.

18 202. GitHub and OpenAI have unjustly utilized access to Licensed Materials hosted on  
19 GitHub. This code is used to create Derivative Works that are licensed to third parties in  
20 exchange for, *inter alia*, compliance with applicable License terms.

21 203. GitHub and OpenAI derive profit or other benefits from removal of attribution,  
22 copyright notices, and license terms from Licensed Materials and reselling it as Output through  
23 Copilot.

24 204. It would be unjust for GitHub and OpenAI to retain those benefits.

25 205. Plaintiffs and the Class have suffered monetary damages as a result of GitHub's  
26 and OpenAI's conduct.



1           212. Plaintiffs and the Class are GitHub users who have accepted GitHub’s Terms of  
2 Service. As a result, Plaintiffs and the Class have formed a contract, the terms of which are set  
3 forth in GitHub’s Terms of Service—including the additional GitHub Copilot Terms from  
4 GitHub Terms for Additional Products and Features.

5           213. Plaintiffs and the Class are GitHub users who have accepted GitHub’s Privacy  
6 Statement. As a result, Plaintiffs and the Class have formed a contract.

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

10           215. Plaintiffs and the Class have performed each of the conditions, covenants, and  
11 obligations imposed on them by the terms of GitHub’s Policies.

12           216. GitHub has substantially and materially breached GitHub’s Policies in the  
13 following ways:

- 14           a. Sharing Plaintiffs’ and the Class’s personal data with unauthorized third parties in  
15 violation of the GitHub Privacy Statement;
- 16           b. Selling and distributing Plaintiffs’ and the Class’s personal data in contravention  
17 of the GitHub Policies;
- 18           c. Use of Plaintiffs’ and the Class’s personal data after the GitHub Privacy Statement  
19 explicitly claims it will be deleted;
- 20           d. Use and distribution of Plaintiffs’ and the Class’s personal data outside the  
21 limitations set forth in the GitHub Privacy Statement.

22           217. Plaintiffs and the Class have suffered monetary damages as a result of GitHub’s  
23 conduct.

24           218. GitHub’s conduct is causing and, unless enjoined and restrained by this Court,  
25 will continue to cause Plaintiffs and the Class great and irreparable injury that cannot fully be  
26 compensated or measured in money.

27  
28





1 they were using, distributing, or selling their PII to unauthorized third parties, namely Copilot  
2 users.

3 226. GitHub and OpenAI also violated the CCPA by failing to provide notice to its  
4 customers of their right to opt-out of the disclosure of their PII to unauthorized third parties,  
5 namely Copilot users.

6 227. GitHub and OpenAI also violated the CCPA by incorporating Plaintiffs' and the  
7 Class's personal information into Copilot with no way to alter or delete. And also with no way to  
8 share that personal data with Plaintiffs or the Class upon request.

9 228. GitHub and OpenAI also violated the CCPA by failing to provide a clear and  
10 conspicuous link entitled "Do Not Sell My Personal Information" to a webpage that enables a  
11 consumer—or a person authorized by a consumer—to opt out of the sale of Plaintiffs' and the  
12 Class's personal data through Copilot.

13 229. By the acts described above, GitHub and OpenAI violated the CCPA by  
14 negligently, carelessly, and recklessly collecting, maintaining, and controlling their customers'  
15 sensitive personal information and by engineering, designing, maintaining, and controlling  
16 systems that exposed their customers' sensitive personal information of which GitHub and  
17 OpenAI had control and possession to the risk of exposure to unauthorized persons, thereby  
18 violating their duty to implement and maintain reasonable security procedures and practices  
19 appropriate to the nature of the information to protect the personal information. GitHub and  
20 OpenAI allowed unauthorized users to view, use, manipulate, exfiltrate, and steal the  
21 nonencrypted and nonredacted personal information of Plaintiffs and other customers, including  
22 their personal and financial information.

23 **COUNT X**  
24 **NEGLIGENCE—NEGLIGENT HANDLING OF PERSONAL DATA**  
25 **Common Law**  
26 **(GitHub and OpenAI)**

27 230. Plaintiffs and the Class hereby repeat and incorporate by reference each preceding  
28 and succeeding paragraph as though fully set forth herein.

1           231.    GitHub and OpenAI owed a duty of reasonable care toward Plaintiffs and the  
2 Class based upon GitHub’s and OpenAI’s relationship to them. This duty is based upon  
3 GitHub’s and OpenAI’s contractual obligations, custom and practice, right to control information  
4 in its possession, exercise of control over the information in its possession, authority to control  
5 the information in its possession, and the commission of affirmative acts that resulted in said  
6 harms and losses. Additionally, this duty is based on the requirements of California Civil Code  
7 section 1714 requiring all “persons,” including GitHub and OpenAI, to act in a reasonable  
8 manner toward others. This duty is also based on the specific statutory duties imposed on  
9 GitHub and OpenAI under California Civil Code sections 1798.100, *et seq.*, as businesses  
10 operating in the State of California that either have annual operating revenue above \$25 million,  
11 collect the personal information of 50,000 or more California residents annually, or derive at least  
12 50 percent of their annual revenue from the sale of personal information of California residents.

13           232.    GitHub and OpenAI breached their duties by negligently, carelessly, and recklessly  
14 collecting, maintaining, and controlling their customers’ sensitive personal information and  
15 engineering, designing, maintaining, and controlling systems—including Copilot—that exposed  
16 and continue to expose their customers’ sensitive personal information of which GitHub and  
17 OpenAI had control and possession to the risk of exposure to unauthorized persons.

18           233.    GitHub and OpenAI also committed per se breaches of said duty by negligently  
19 violating the dictates of California Civil Code sections 1798.82, *et seq.*, and 1798.100, *et seq.*, and  
20 the provisions of the California Constitution enshrining the right to privacy, by failing to inform  
21 Plaintiffs and the Class of the access to their sensitive personal information by unauthorized  
22 persons expeditiously and without delay and failing to adequately safeguard this information from  
23 unauthorized access even after GitHub and OpenAI became aware of multiple instances of  
24 release of this information by Copilot. The provisions of the California Civil Code and the  
25 California Constitution that GitHub and OpenAI violated were enacted to protect the class of  
26 Plaintiffs here involved from the type of injury here incurred, namely their right to privacy and  
27 the protection of their personal data. Plaintiffs and the Class were within the class of persons and  
28

1 consumers who were intended to be protected by California Civil Code sections 1798.82, *et seq.*,  
2 and 1798.100, *et seq.*

3 234. As a direct consequence of the actions described herein, and the breaches of  
4 duties indicated thereby, unauthorized users gained access to, exfiltrated, stole, and gained  
5 disclosure of the sensitive personal information of Plaintiffs and the Class, causing them harms  
6 and losses including but not limited to economic loss, the loss of control over the use of their  
7 identity, harm to their constitutional right to privacy, lost time dedicated to cure harm to their  
8 privacy, the need for future expenses and time dedicated to the recovery and protection of further  
9 loss, and privacy injuries associated with having their sensitive personal and financial information  
10 disclosed.

11 **COUNT XI**  
12 **CIVIL CONSPIRACY**  
13 **Common Law**  
**(Against All Defendants)**

14 235. Plaintiffs and the Class hereby repeat and incorporate by reference each preceding  
15 and succeeding paragraph as though fully set forth herein.

16 236. On information and belief, Microsoft, GitHub, OpenAI, and the Individual  
17 Defendants have worked together to create Copilot. In creating Copilot, Defendants willfully  
18 avoided determining whether and how Copilot's training and Output may violate the rights of  
19 Plaintiffs and the Class and other stakeholders. This is because Defendants understood that  
20 through Copilot they would be engaging in a variety of unlawful conduct. Defendants conduct  
21 resulted in violations of Plaintiffs' and the Class's rights as set forth herein.

22 237. On information and belief, OpenAI derives a financial or other valuable benefit  
23 from the sale of Copilot. In exchange, OpenAI provided Microsoft an exclusive license to use its  
24 GPT-3 language model.

25 238. On information and belief, Microsoft derives a financial benefit from sales of  
26 Copilot through payments or other form of compensation in exchange for GitHub's and  
27 OpenAI's use of Azure to run Copilot.  
28



- 1 a) Judgment in favor of Plaintiffs and the Class and against Defendants;
- 2 b) Permanent injunctive relief, including but not limited to making changes to its
- 3 Copilot product to ensure that all applicable information set forth in 17 U.S.C. §
- 4 1203(b)(1) is included in along with any Output including associated code;
- 5 c) An order of costs and allowable attorney's fees pursuant to 17 U.S.C. §
- 6 1203(b)(4)–(5);
- 7 d) An award of statutory damages pursuant to 17 U.S.C. § 1203(b)(3) and 17 U.S.C. §
- 8 1203(c)(3),<sup>41</sup> or, in the alternative, an award of actual damages and any additional
- 9 profits pursuant to 17 U.S.C. § 1203(c)(2) (including tripling damages pursuant to
- 10 17 U.S.C. § 1203(c)(4) if applicable);
- 11 e) An award of damages for harms resulting from Defendants' breach of Licenses;
- 12 f) An award of damages, including punitive damages, for harms resulting from
- 13 Defendants' tortious interference in Plaintiffs' and the Class's prospective
- 14 contractual relations;
- 15 g) An award of damages for harms resulting from Defendants' false designation of
- 16 the origin of Copilot's Output;
- 17 h) An award of damages in the amount Defendants have been unjustly enriched
- 18 through their conduct as alleged herein as well as punitive damages in connection
- 19 with this conduct;
- 20 i) An award of damages, including punitive damages, for harms resulting from
- 21 Defendants acts of unfair competition;
- 22 j) Statutory damages and any other relief this Court deems proper for Defendants

---

23 <sup>41</sup> Plaintiffs estimate that statutory damages for Defendants' direct violations of DMCA Section

24 1202 alone will exceed \$9,000,000,000. That figure represents minimum statutory damages

25 (\$2,500) incurred three times for each of the 1.2 million Copilot users Microsoft reported in June

26 2022. Each time Copilot provides an unlawful Output it violates Section 1202 three times

27 (distributing the Licensed Materials without: (1) attribution, (2) copyright notice, and (3) License

28 Terms). So, if each user receives just one Output that violates Section 1202 throughout their time

using Copilot (up to fifteen months for the earliest adopters), then GitHub and OpenAI have

violated the DMCA 3,600,000 times. At minimum statutory damages of \$2500 per violation, that

translates to \$9,000,000,000.

1 violation of the CCPA;

2 k) An award of damages for harms resulting from GitHub’s breach of the GitHub  
3 Policies; and

4 l) An award of damages, including punitive damages, for harms resulting from  
5 Defendants’ negligent handling of Plaintiffs’ and the Class’s personal data.

6 244. Injunctive relief sufficient to alleviate and stop Defendants’ unlawful conduct  
7 alleged herein.

8 245. Plaintiffs and the Class are entitled to prejudgment and post-judgment interest on  
9 the damages awarded them, and that such interest be awarded at the highest legal rate from and  
10 after the date this class action complaint is first served on Defendants;

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

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

16 **X. JURY TRIAL DEMANDED**

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

19 ///

20 ///

21 ///

22

23

24

25

26

27

28

1 Dated: November 3, 2022

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

2  
3  
4 Joseph R. Saveri (State Bar No. 130064)  
5 Cadio Zirpoli (State Bar No. 179108)  
6 Travis Manfredi (State Bar No. 281779)  
7 **JOSEPH SAVERI LAW FIRM, LLP**  
8 601 California Street, Suite 1000  
9 San Francisco, California 94108  
10 Telephone: (415) 500-6800  
11 Facsimile: (415) 395-9940  
12 Email: jsaveri@saverilawfirm.com  
13 czirpoli@saverilawfirm.com  
14 tmanfredi@saverilawfirm.com

15  
16 Matthew Butterick (State Bar No. 250953)  
17 1920 Hillhurst Avenue, #406  
18 Los Angeles, CA 90027  
19 Telephone: (323) 968-2632  
20 Facsimile: (415) 395-9940  
21 Email: mb@buttericklaw.com  
22 *Counsel for Plaintiffs and the Proposed Class*  
23  
24  
25  
26  
27  
28