

1 Joseph R. Saveri (State Bar No. 130064)
 Steven N. Williams (State Bar No. 175489)
 2 Cadio Zirpoli (State Bar No. 179108)
 Christopher K.L. Young (State Bar No. 318371)
 3 Louis A. Kessler (State Bar No. 243703)
 Elissa A. Buchanan (State Bar No. 249996)
 4 Travis Manfredi (State Bar No. 281779)
 5 **JOSEPH SAVERI LAW FIRM, LLP**
 601 California Street, Suite 1000
 San Francisco, California 94108
 Telephone: (415) 500-6800
 Facsimile: (415) 395-9940
 8 Email: jsaveri@saverilawfirm.com
 swilliams@saverilawfirm.com
 9 czirpoli@saverilawfirm.com
 cyoung@saverilawfirm.com
 10 lkessler@saverilawfirm.com
 eabuchanan@saverilawfirm.com
 11 tmanfredi@saverilawfirm.com
 12

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

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

18 **UNITED STATES DISTRICT COURT**
 19 **NORTHERN DISTRICT OF CALIFORNIA**
 20 **OAKLAND DIVISION**

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

22 Individual and Representative Plaintiffs,

23 v.

24 GITHUB, INC., a Delaware corporation;
 MICROSOFT CORPORATION, a Washington
 25 corporation; OPENAI, INC., a Delaware nonprofit
 corporation; OPENAI, L.P., a Delaware limited
 26 partnership; OPENAI OPCO, L.L.C., a Delaware
 limited liability company; OPENAI GP, L.L.C., a
 27 Delaware limited liability company; OPENAI
 28 STARTUP FUND GP I, L.L.C., a Delaware limited

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

FIRST AMENDED COMPLAINT

CLASS ACTION

DEMAND FOR JURY TRIAL

1 liability company; OPENAI STARTUP FUND I, L.P., a
2 Delaware limited partnership; OPENAI STARTUP
3 FUND MANAGEMENT, LLC, a Delaware limited
4 liability company,

Defendants.

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

TABLE OF CONTENTS

1
2
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 PIRACY1

II. JURISDICTION AND VENUE..... 4

III. INTRADISTRICT ASSIGNMENT 4

IV. PARTIES..... 4

 A. Plaintiffs 4

 B. Defendants 6

V. AGENTS AND CO-CONSPIRATORS 8

VI. CLASS ALLEGATIONS 9

 A. Class Definitions..... 9

 B. Numerosity.....10

 C. Typicality.....10

 D. Commonality & Predominance.....10

 1. DMCA Violations10

 2. Contract-Related Conduct 11

 3. Unlawful-Competition Conduct 11

 4. Injunctive Relief..... 11

 5. Defenses 11

 E. Adequacy.....12

 F. Other Class Considerations12

VII. FACTUAL ALLEGATIONS12

 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 17

 D. Codex and Copilot Were Trained on Copyrighted Materials Offered
Under Licenses..... 20

 E. Copilot Was Launched Despite Its Propensity for Producing Unlawful
Outputs21

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

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

- 1. Example: Copilot Outputs the Code of Doe 2 Essentially Verbatim 24
- 2. Example: Copilot Outputs the Code of Doe 1 in Modified Format 26
- 3. Example: Copilot Outputs the Code of Doe 5 In Modified Format 30
- 4. Example: Copilot Outputs Code of Doe 5 Essentially Verbatim33
- G. Codex and Copilot Were Designed to Withhold Attribution, Copyright Notices, and License Terms from Their Users37
- H. Open-Source Licenses Began to Appear in the Early 1990s39
- I. Microsoft Has a History of Flouting Open-Source License Requirements.....41
- J. GitHub Was Designed to Cater to Open-Source Projects 43
- K. OpenAI Is Intertwined with Microsoft and GitHub..... 45
- L. Conclusion of Factual Allegations 47
- VIII. CLAIMS FOR RELIEF..... 48
- IX. DEMAND FOR JUDGMENT 64
- X. JURY TRIAL DEMANDED 66

1 Plaintiffs J. Doe 1, J. Doe 2, J. Doe 3, J. Doe 4 and J. Doe 5 (“Plaintiffs”), on behalf of
2 themselves and all others similarly situated, bring this Class Action Complaint (the “Complaint”)
3 against Defendants GitHub, Inc.; Microsoft Corporation; OpenAI, Inc.; OpenAI, L.P.; OpenAI
4 OpCo, L.L.C; OpenAI GP, L.L.C.; OpenAI Startup Fund GP I, L.L.C.; OpenAI Startup Fund I,
5 L.P.; and OpenAI Startup Fund Management, LLC¹ for violation of the Digital Millennium
6 Copyright Act, 17 U.S.C. §§ 1201–1205 (the “DMCA”); breach of contract regarding the
7 Suggested Licenses, breach of contract regarding GitHub’s policies including its terms of service;
8 tortious interference with prospective economic relations; California’s Unfair Competition law,
9 Cal. Bus. & Prof. Code section 17200, *et seq.*; common law unfair competition; negligence, and
10 unjust enrichment.

11 I. OVERVIEW: A BRAVE NEW WORLD 12 OF SOFTWARE PIRACY

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

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

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

24 ¹ GitHub, Inc. is referred to as “GitHub.” Microsoft Corporation is referred to as “Microsoft.”
25 OpenAI, Inc.; OpenAI, L.P.; OpenAI OpCo, L.L.C.; OpenAI GP, L.L.C.; OpenAI Startup Fund
26 GP I, L.L.C.; OpenAI Startup Fund I, L.P.; and OpenAI Startup Fund Management, LLC are
27 referred to collectively herein as “OpenAI.” Collectively, GitHub, Inc., Microsoft Corporation,
28 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 open-source source code on the website github.com. Over the next 10 years, GitHub, based on
2 these representations succeeded wildly, attracting nearly 25 million developers.

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

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

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

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

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

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

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

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

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

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

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

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

23
24
25 ² "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 ³ 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.

1 **II. JURISDICTION AND VENUE**

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

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

15 **III. INTRADISTRICT ASSIGNMENT**

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

22 **IV. PARTIES**

23 **A. PLAINTIFFS**

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

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

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

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

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

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

B. DEFENDANTS

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

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

26. Defendant OpenAI, Inc. is a Delaware nonprofit corporation with its principal place of business located at 3180 18th Street, San Francisco, CA 94110. OpenAI, Inc. is a party to the unlawful conduct alleged herein. It—along with OpenAI, L.P.—programed, trained, and maintains Codex, which infringes all the same rights at Copilot and is also an integral piece of Copilot. Copilot requires Codex to function. OpenAI, Inc. is a party to the unlawful conduct alleged herein. OpenAI, Inc. founded, owns, and exercises control over all the other OpenAI entities, including those set forth in Paragraphs 27–32.

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

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

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

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

27 31. Defendant OpenAI Startup Fund GP I, L.L.C. ("OpenAI Startup Fund GP I") is
28 a Delaware limited liability company with its principal place of business located at 3180 18th

1 Street, San Francisco, CA 94110. OpenAI Startup Fund GP I is the general partner of OpenAI
2 Startup Fund I. OpenAI Startup Fund GP I manages and operates the day-to-day business and
3 affairs of OpenAI Startup Fund I. OpenAI Startup Fund GP I is liable for the debts, liabilities and
4 obligations of OpenAI Startup Fund I, including litigation and judgments. OpenAI Startup Fund
5 GP I was aware of the unlawful conduct alleged herein and exercised control over OpenAI, L.P.
6 throughout the Class Period. OpenAI Startup Fund GP I is a party to the unlawful conduct
7 alleged herein. Sam Altman, co-founder, CEO, and Board member of OpenAI, Inc. is the
8 Manager of OpenAI Startup Fund GP I. OpenAI Startup Fund GP I is the General Partner of
9 OpenAI Startup Fund I, L.P.

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

16 V. AGENTS AND CO-CONSPIRATORS

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

21 34. The Defendants’ agents operated under the explicit and apparent authority of
22 their principals.

23 35. Each Defendant, and its subsidiaries, affiliates and agents operated as a single
24 unified entity.

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

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

VI. CLASS ALLEGATIONS

A. Class Definitions

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

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

All persons or entities domiciled in the United States that, (1) owned an interest in at least one US copyright in any work; (2) offered that work under one of GitHub’s Suggested Licenses⁴; and (3) stored Licensed Materials in any public GitHub repositories at any time between January 1, 2015 and the present (the “Class Period”).

“Damages Class” under Rule 23(b)(3):

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

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

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

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

B. Numerosity

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

C. Typicality

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

D. Commonality & Predominance

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

1. DMCA Violations

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

1 **2. Contract-Related Conduct**

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

11 **3. Unlawful-Competition Conduct**

- 12 • Whether Defendants passed-off the Licensed Materials as its own creation
13 and/or Codex or Copilot’s creation.
14 • Whether Defendants were unjustly enriched by the unlawful conduct alleged
15 herein.
16 • Whether Defendants’ conduct alleged herein constitutes Unfair Competition
17 under California Business and Professions Code section 17200 *et seq.*
18 • Whether Defendants’ conduct alleged herein constitutes unfair competition
19 under the common law.

20 **4. Injunctive Relief**

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

23 **5. Defenses**

- 24 • Whether any affirmative defense excuses Defendants’ conduct.
25 • Whether any statutes of limitation limit Plaintiffs’ and the Class’s potential for
26 recovery.
27 • Whether any applicable statutes of limitation should be tolled as a result of
28 Defendants’ fraudulent concealment of their unlawful conduct.

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

3 **E. Adequacy**

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

10 **F. Other Class Considerations**

11 44. Defendants have acted on grounds generally applicable to the Class, thereby
12 making final injunctive relief appropriate with respect to the Class as a whole.

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

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

20 **VII. FACTUAL ALLEGATIONS**

21 **A. Introduction**

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

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

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

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

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

16 **B. Codex Outputs Copyrighted Materials Without Following the Terms of the**
17 **Applicable Licenses**

18 52. Below is an explanation of how Codex functions. When Codex is prompted with:
19 function isEven(n) {
20 it assumes this is the beginning of a function written in the JavaScript language that will test
21 whether a number is even.

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

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

1 53. 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:⁷

```
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     // → ??
```

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

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

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

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

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

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

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

26 ⁸ <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 61. Here is the exercise from *Eloquent JavaScript*:

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

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

14 63. Haverbeke provides the following solution to the function discussed above:

```
15     function isEven(n) {  
16         if (n == 0) return true;  
17         else if (n == 1) return false;  
18         else if (n < 0) return isEven(-n);  
19         else return isEven(n - 2);  
20     }  
21     console.log(isEven(50));  
22     // → true  
23     console.log(isEven(75));  
24     // → false  
25     console.log(isEven(-1));  
26     // → false
```

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

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

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

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

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

15 69. Ultimately, Codex derives its value primarily from its ability to locate and output
16 potentially useful Licensed Materials. And from its obfuscation of any rights associated with
17 those materials.

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

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

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

26
27
28 ⁹ See Appendix A for full text of the MIT License.

1 72. When we give Copilot the same prompt discussed above in Paragraph 52,
2 “function isEven(n) {”, it interprets the prompt as the beginning of a function written in
3 the JavaScript language that will test whether a number is even, same as Codex.

4 73. However, the Output of Copilot in response to the prompt is different than Codex,
5 namely:

```
6     function isEven(n) {
7         return n % 2 === 0;
8     }
```

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

12 75. Copilot’s Output, like Codex’s, is derived from existing code. Namely, sample
13 code that appears in the online book *Mastering JS*, written by Valeri Karpov.¹⁰ Like Codex’s
14 Output, Copilot’s is also based upon copyrighted educational material. *Mastering JS* is a set of
15 educational exercises for programmers. Like *Eloquent JavaScript*, there are many copies of
16 Karpov’s exercise stored in public repositories on GitHub. Programmers working through
17 *Mastering JS* store their answers there.

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

```
21     function isPrime(n) {
22         if (n < 2) {
23             return false;
24         }
25         for (let i = 2; i < n; i++) {
26             if (n % i === 0) {
27                 return false;
28             }
29         }
30         return true;
31     }
```

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

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

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

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

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

22 80. The material in Curinga’s book is made available under the GNU Free
 23 Documentation License. Although this is not one of the Suggested Licenses, it contains similar
 24 attribution provisions, namely that “You may copy and distribute the Document in any medium,
 25 either commercially or noncommercially, provided that this License, the copyright notices, and
 26
 27

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

1 the license notice saying this License applies to the Document are reproduced in all copies, and
2 that you add no other conditions whatsoever to those of this License.”¹²

3 81. As with Codex, Copilot does not provide the end user any attribution of the
4 original author of the code, nor anything about their license requirements. There is no way for the
5 Copilot user to know that they must provide attribution, copyright notice, nor a copy of the
6 license’s text. And with regard to the GNU Free Documentation License, Copilot users would
7 not be aware that they are limited in what conditions they can place on the use of derivative works
8 they make using this copyrighted code. Had the Copilot user found this code in a public GitHub
9 repository or a copy of the book it was originally published in, they would find the GNU Free
10 Documentation License at the same time and be aware of its terms. Copilot finds that code for the
11 user but excises the license terms, copyright notice, and attribution. This practice allows its users
12 to assume that the code can be used without restriction. It cannot.

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

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

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

26
27
28 ¹² <https://matt.curinga.com/think-js/#gnu-free-documentation-license>.

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

4 85. The words “study” and “training” and “learning” in connection with AI describe
5 algorithmic processes that are not analogous to human reasoning. AI models cannot “learn” as
6 humans do, nor can it “understand” semantics and context the way humans do. Rather, it detects
7 statistically significant patterns in its training data and provides Output derived from its training
8 data when statistically appropriate. A “brute force” approach like this would not be efficient nor
9 even possible for humans. A human could not memorize, statistically analyze, and easily access
10 thousands of gigabytes of existing code, a task now possible for powerful computers like those
11 that make up Microsoft’s Azure cloud platform. To accomplish the same task, a human may
12 search for Licensed Materials that serve their purpose if they believe such materials exist. And if
13 that human finds such materials, they will probably abide by its License Terms rather than risk
14 infringing its owners’ rights. At the very least, if they incorporate those Licensed Materials into
15 their own project without following its terms they will be doing so knowingly.

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

17 86. GitHub and OpenAI have not provided much detail regarding what data Codex
18 and OpenAI were trained on. Plaintiffs know for certain from GitHub and OpenAI’s statements,
19 that both systems were trained on publicly available GitHub repositories, with Copilot having
20 been trained on all available public GitHub repositories.

21 87. According to OpenAI, Codex was trained on “billions of lines of source code from
22 publicly available sources, including code in public GitHub repositories.” Similarly, GitHub has
23 described¹³ Copilot’s training material as “billions of lines of public code.” GitHub researcher
24 Eddie Aftandilian confirmed in a recent podcast¹⁴ that Copilot is “train[ed] on public repos on
25 GitHub.”

26
27 ¹³ <https://github.blog/2021-06-30-github-copilot-research-recitation/>.

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

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

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

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

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

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

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

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

1 91. Beyond the examples above, Copilot regularly Output’s verbatim copies of
2 Licensed Materials. For example, Copilot reproduced verbatim well-known code from the game
3 Quake III, use of which is governed by one of the Suggested Licenses—GPL-2.¹⁷

4 92. Copilot also reproduced code that had been released under a license that allowed
5 its use only for free games and required attribution by including a copy of the license. Copilot did
6 not mention nor include the underlying license when providing a copy of this code as Output.¹⁸

7 93. Texas A&M computer-science professor Tim Davis has provided numerous
8 examples of Copilot reproducing code belonging to him without its license or attribution.¹⁹

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

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

23
24 ¹⁷ <https://twitter.com/stefankarpinski/status/1410971061181681674/>.

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

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

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

1 96. Furthermore, the Suggested Licenses impose attribution obligations not only
2 when Licensed Materials have been used verbatim, but also when Licensed Materials have been
3 modified or adapted. Though Output from Copilot is often a verbatim copy, even more often it is
4 a modification: for instance, a near-identical copy that contains only semantically insignificant
5 variations of the original Licensed Materials, or a modified copy that recreates the same
6 algorithm. Whenever Copilot outputs Licensed Materials in a manner that qualifies as a
7 modification, the attribution requirements of the Suggested Licenses still apply. Copilot’s failure
8 to provide the attributions for outputs that are modifications of Licensed Materials represents
9 another enormous set of license breaches and DMCA violations.

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

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

14 98. Described below are some specific examples of Copilot’s unlawful behavior using
15 Licensed Materials owned by the named Plaintiffs. These examples were emitted by Copilot after
16 prompting Copilot.

17 99. In the examples below, original code is shaded gray, prompts to Copilot are shaded
18 orange, and outputs from Copilot are shaded light blue.

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

20 100. The first example demonstrates Copilot suggesting an essentially verbatim copy of
21 code written by Doe 2.

22 101. [REDACTED]

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

24 [REDACTED]

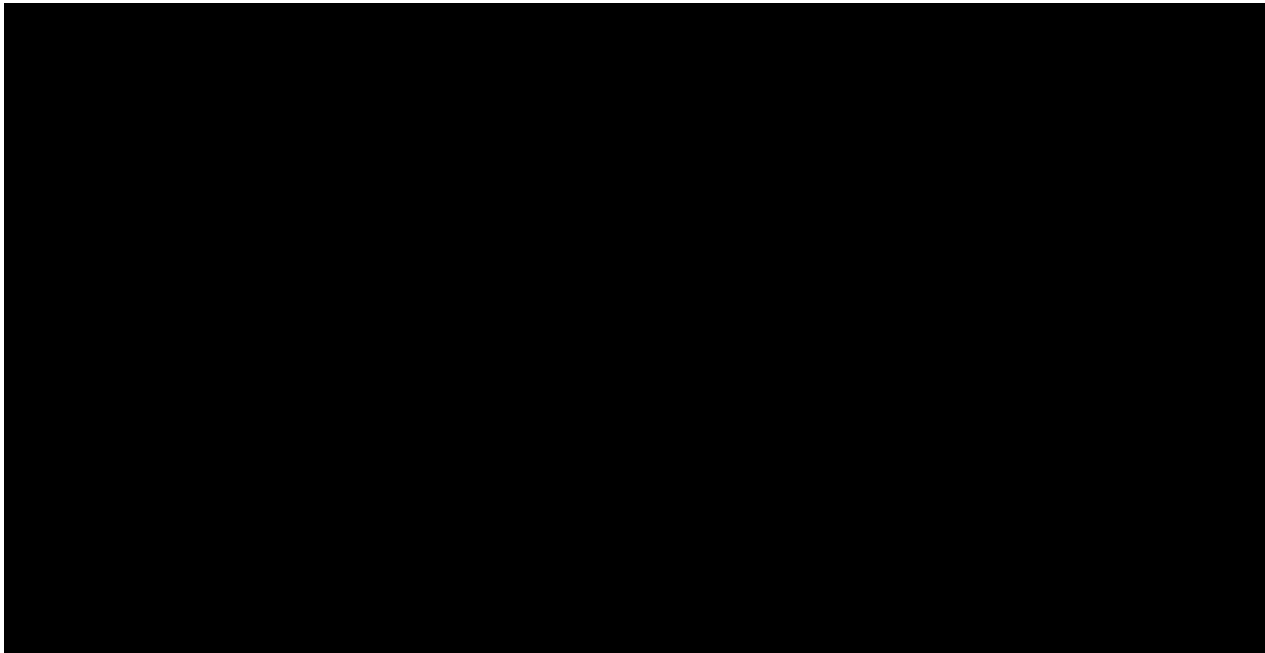
25 [REDACTED]

26 [REDACTED]:

27 [REDACTED]

28 [REDACTED]

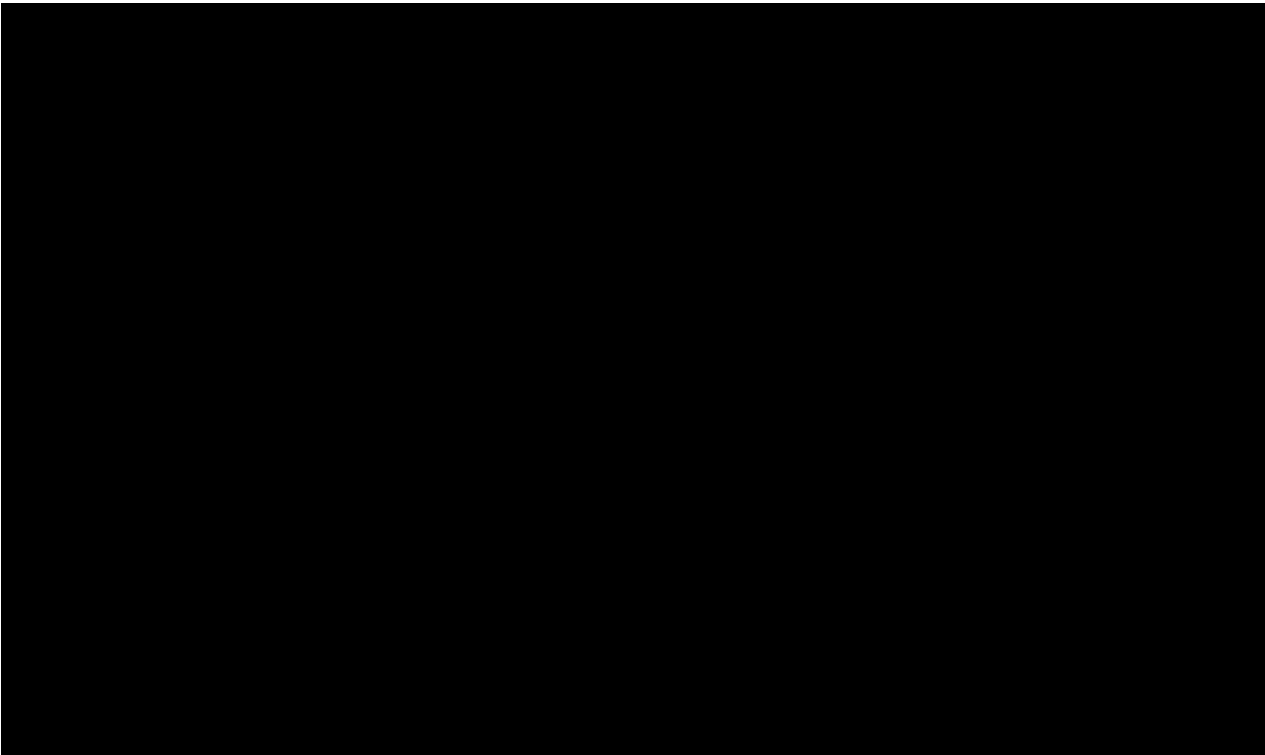
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28



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



Copilot suggests the following:



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

2 [REDACTED]

3 [REDACTED]

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

7 [REDACTED]

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

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

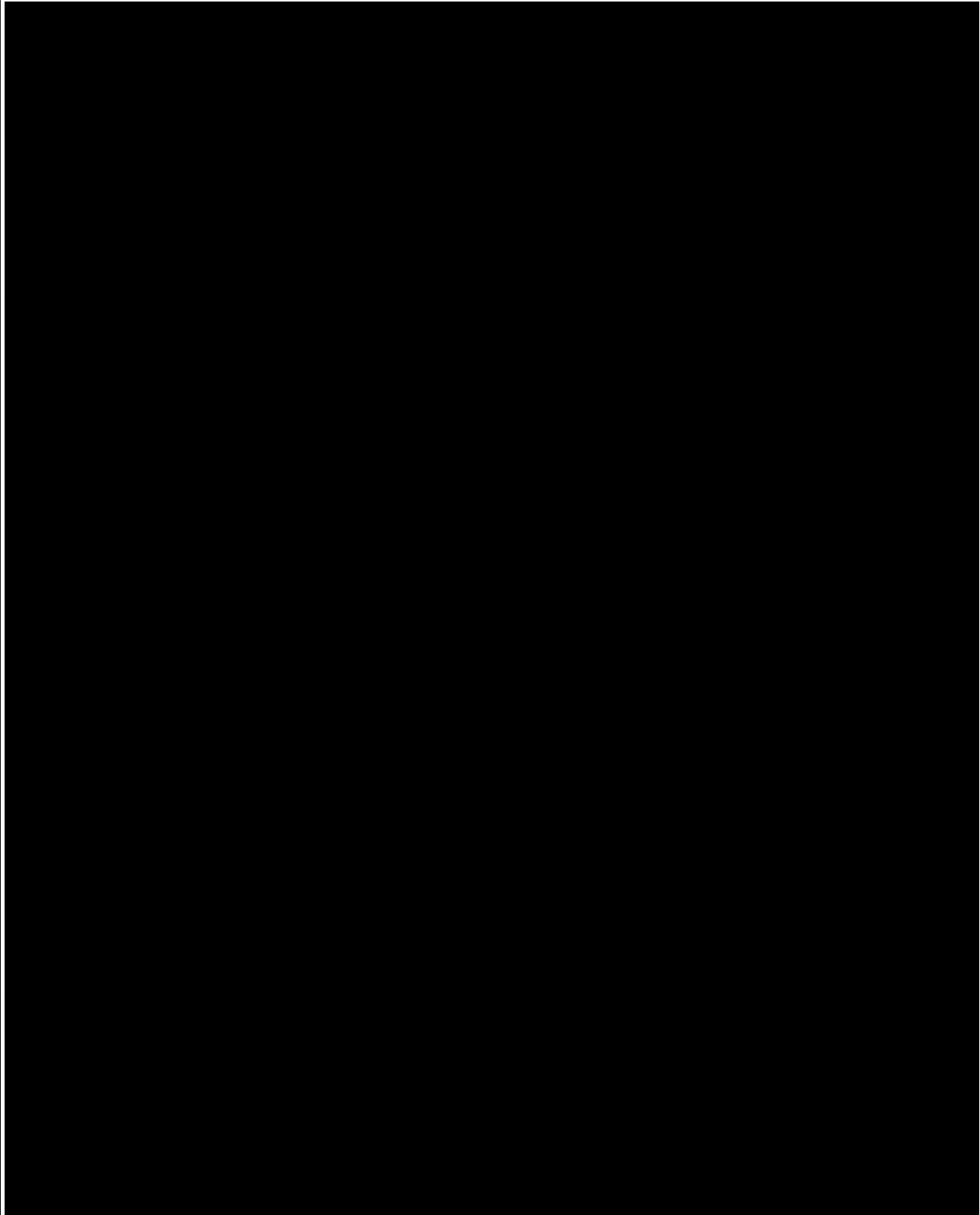
13 105. The second example demonstrates Copilot suggesting a modified copy of code
14 written by Doe 1. To protect Doe 1’s identity, the paragraphs describing the code will be redacted.

15 106. [REDACTED]
16 subject to the MIT License. [REDACTED]

17 [REDACTED]

18 [REDACTED]
19 [REDACTED]
20 [REDACTED]
21 [REDACTED]
22 [REDACTED]
23 [REDACTED]
24 [REDACTED]
25 [REDACTED]
26 [REDACTED]
27 [REDACTED]
28 [REDACTED]

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

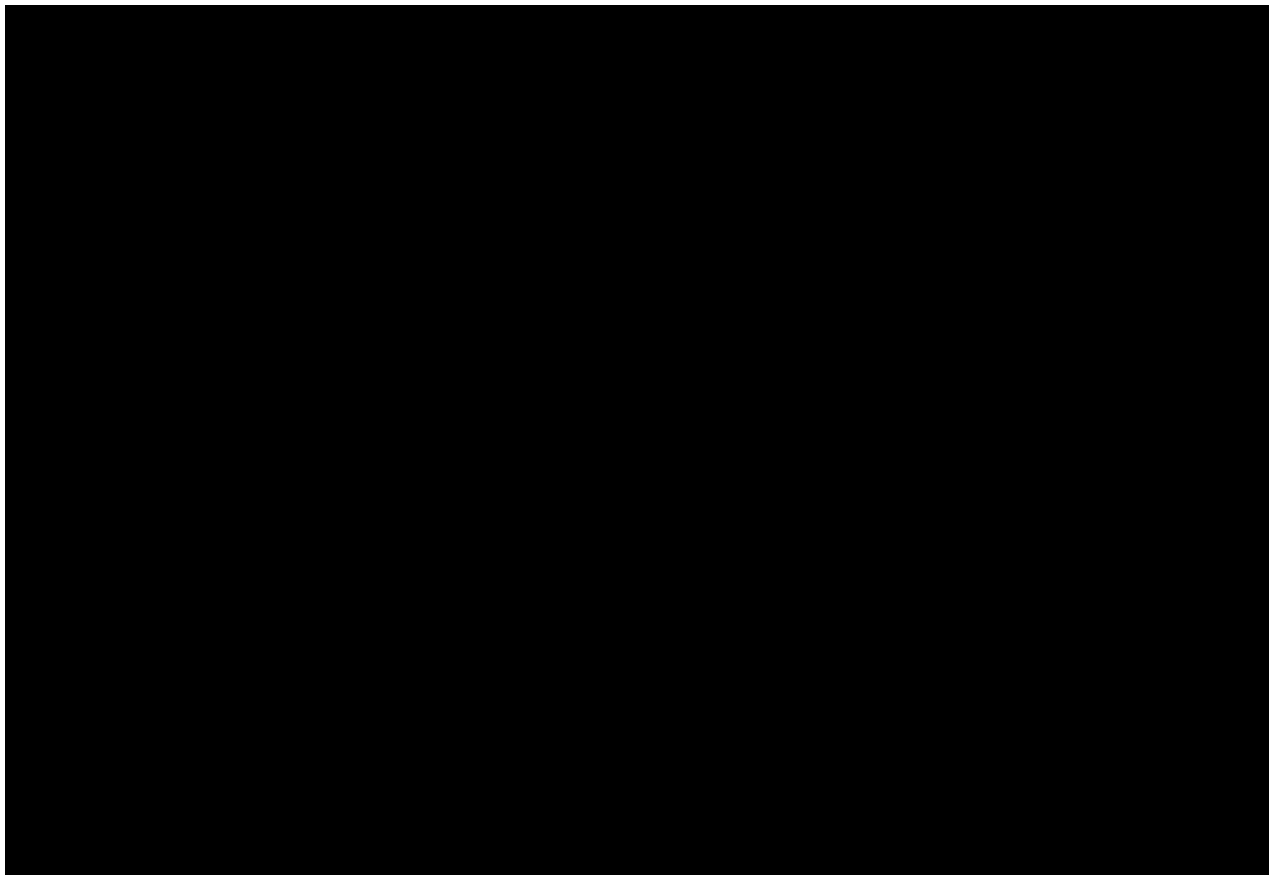


107. When Copilot is prompted with [REDACTED] [REDACTED] [REDACTED]

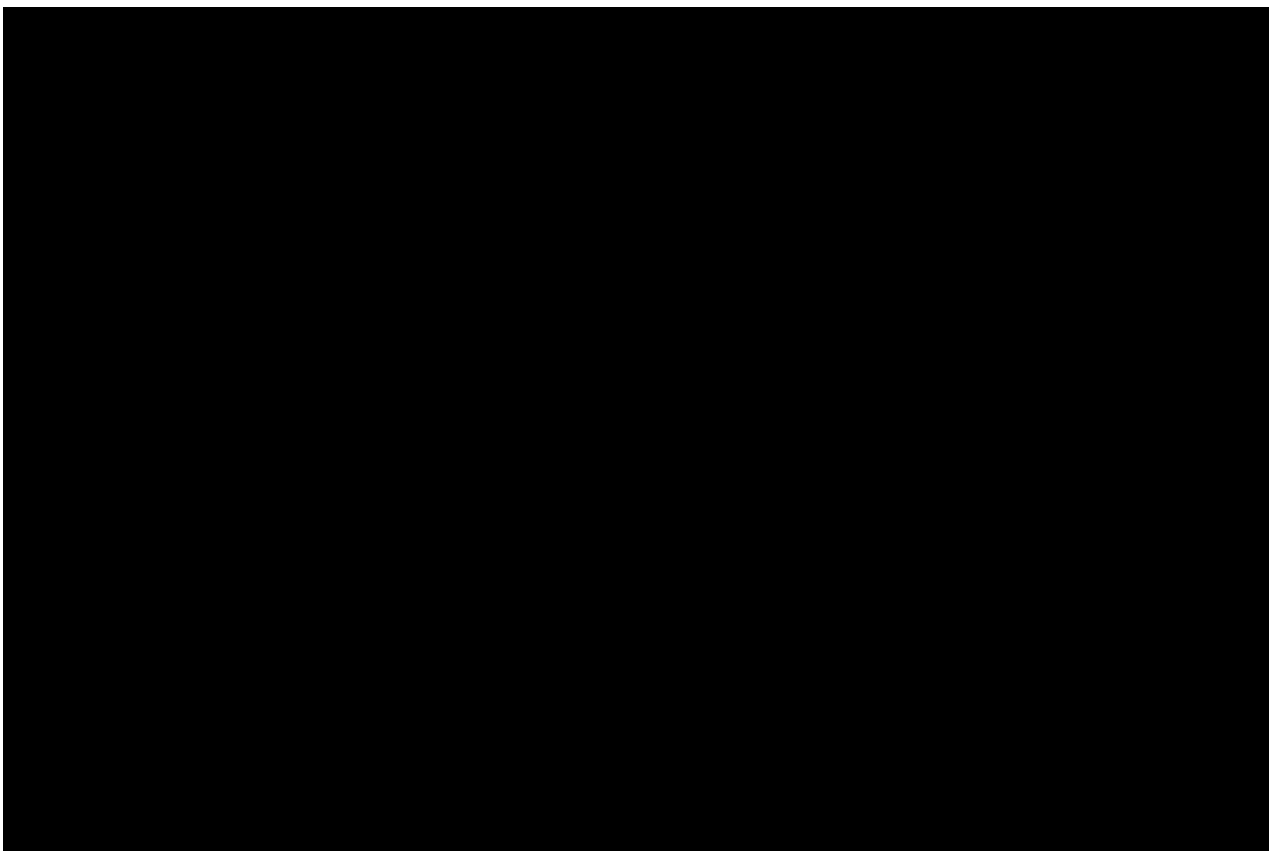
[REDACTED] [REDACTED] [REDACTED]

[REDACTED]

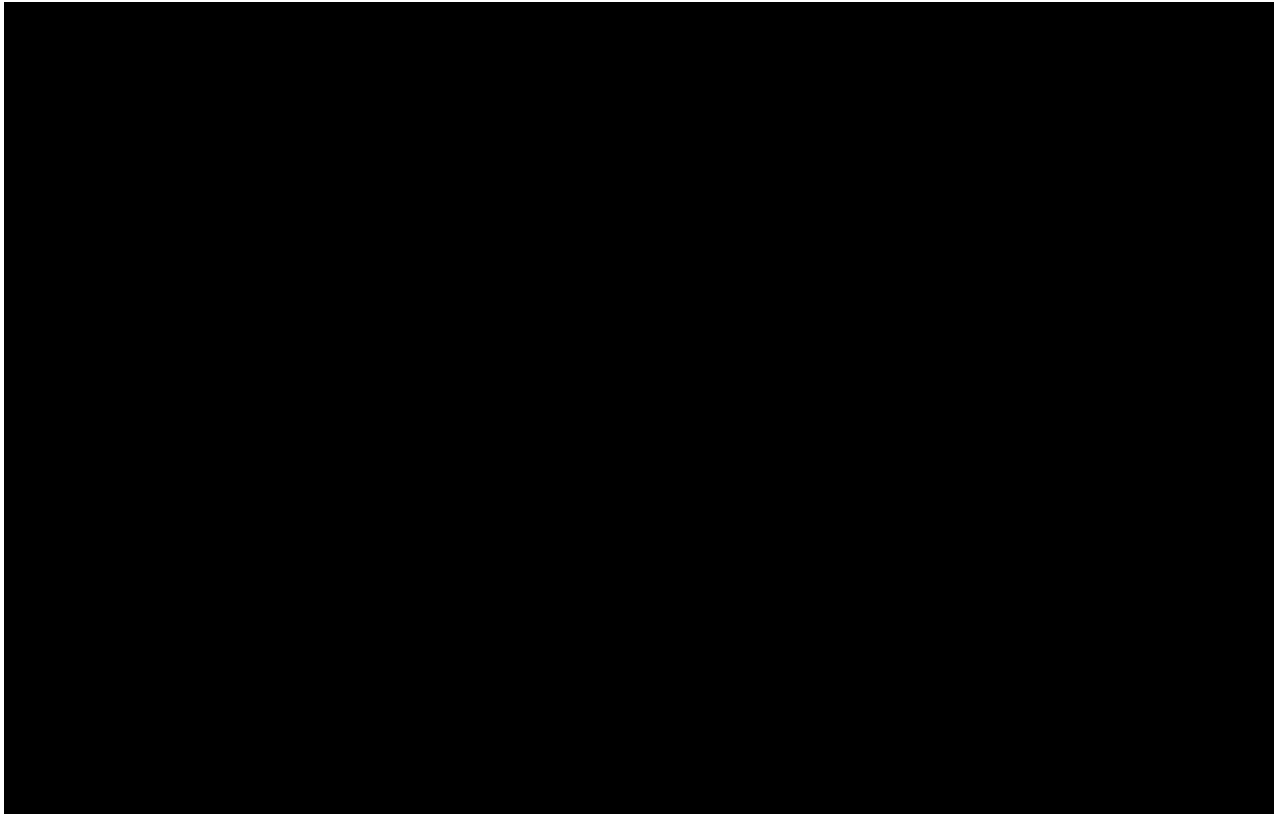
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28



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



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



108. [REDACTED]

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

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

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

means the same thing as this Copilot-suggested code:

[REDACTED]

1 110. As is apparent from a cursory glance of this example, the variations between
2 Copilot’s emitted output and Doe 1’s source code are cosmetic and the code is functionally
3 equivalent; it follows that Copilot’s output is a copy of Doe 1’s code.

4 111. That said, Copilot also introduces mistakes into the code. For instance, [REDACTED]
5 [REDACTED] [REDACTED] [REDACTED] [REDACTED]
6 [REDACTED]

7 112. Still, because Copilot is reproducing Doe 1’s algorithm in modified format, and the
8 obligations in Doe 1’s license (the MIT License) carry with the code even if the underlying code
9 is modified, the Copilot suggestion needs to follow the requirements of Doe 1’s license for that
10 code, including providing attribution. It does not. Copilot also did not reproduce Doe 1’s license.

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

12 113. The third example demonstrates Copilot suggesting multiple modified copies of
13 code written by Doe 5 in response to a sequence of prompts, which is a common way of using
14 Copilot. To protect Doe 5’s identity, the paragraphs describing the code will be redacted.

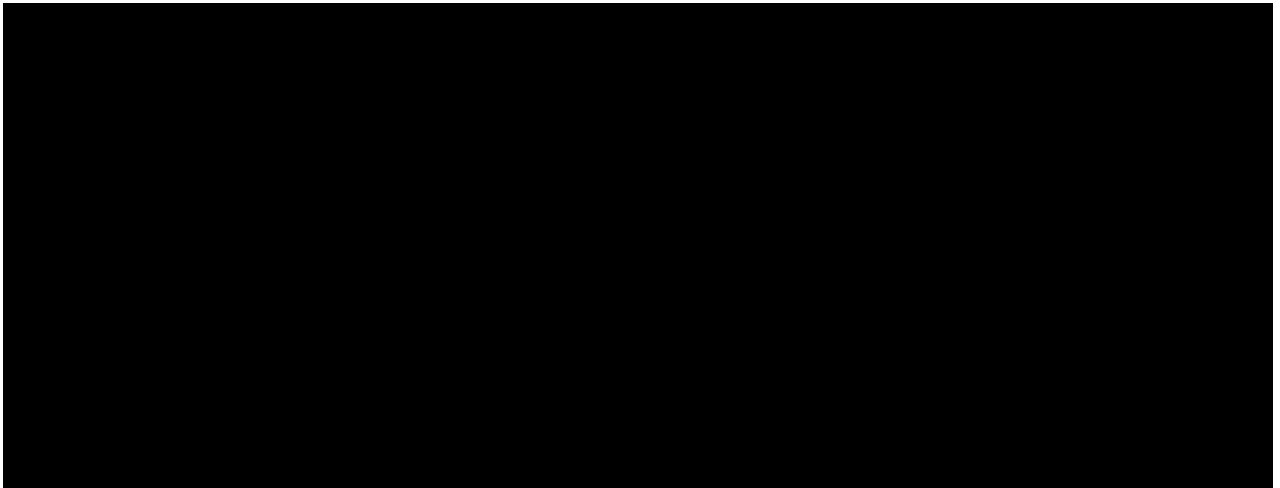
15 114. [REDACTED]
16 [REDACTED] subject to the MIT License. [REDACTED]
17 [REDACTED] The relevant code from the original source file is shown below:

18 [REDACTED]
19 [REDACTED]
20 [REDACTED]
21 [REDACTED]
22 [REDACTED]
23 [REDACTED]
24 [REDACTED]
25 [REDACTED]
26 [REDACTED]
27 [REDACTED]
28 [REDACTED]

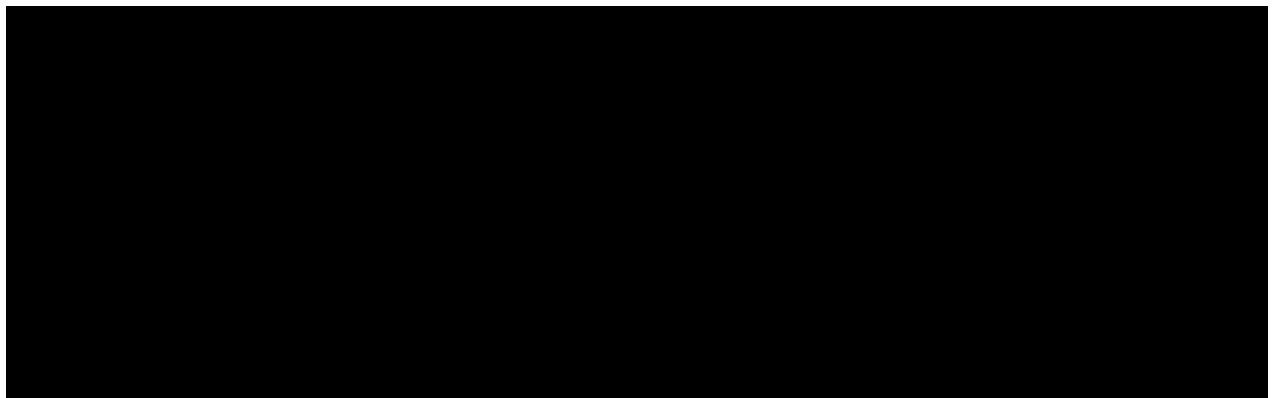
1 115. When Copilot is prompted the first section of Doe 5's code, comprising the first
2 complete test and the name of the second:



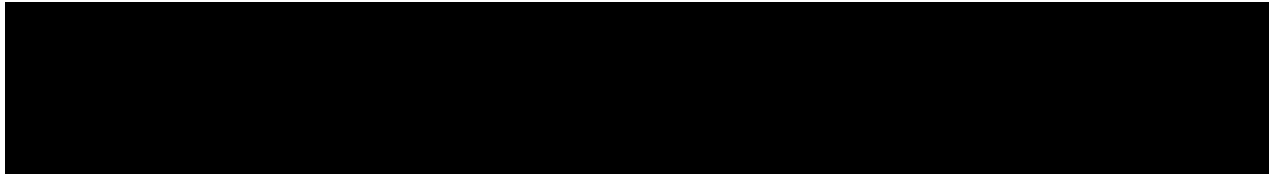
9 116. The first suggestion from Copilot offers to complete the prompt with a verbatim
10 copy of Doe 5's original code, [REDACTED] [REDACTED] [REDACTED]
11 [REDACTED] (a variation that does not affect how the code works):



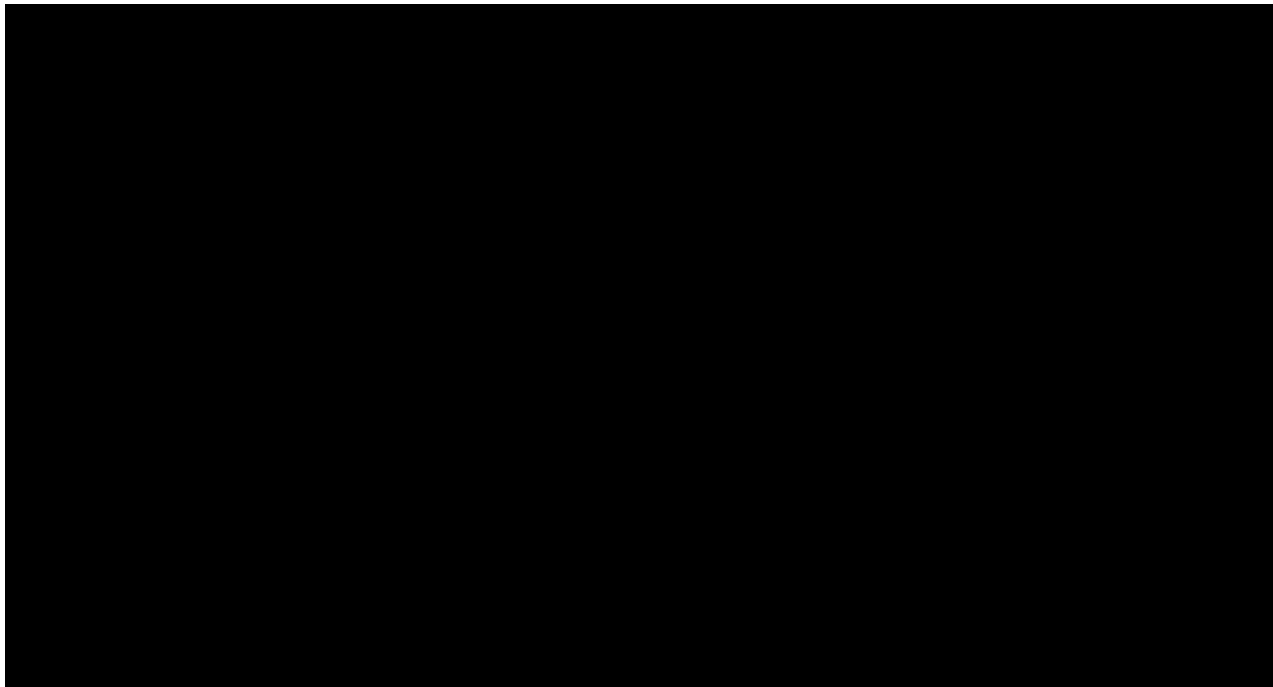
20 117. Next, if the name of the third test is appended, the next prompt to Copilot looks
21 like this:



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



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



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

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

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

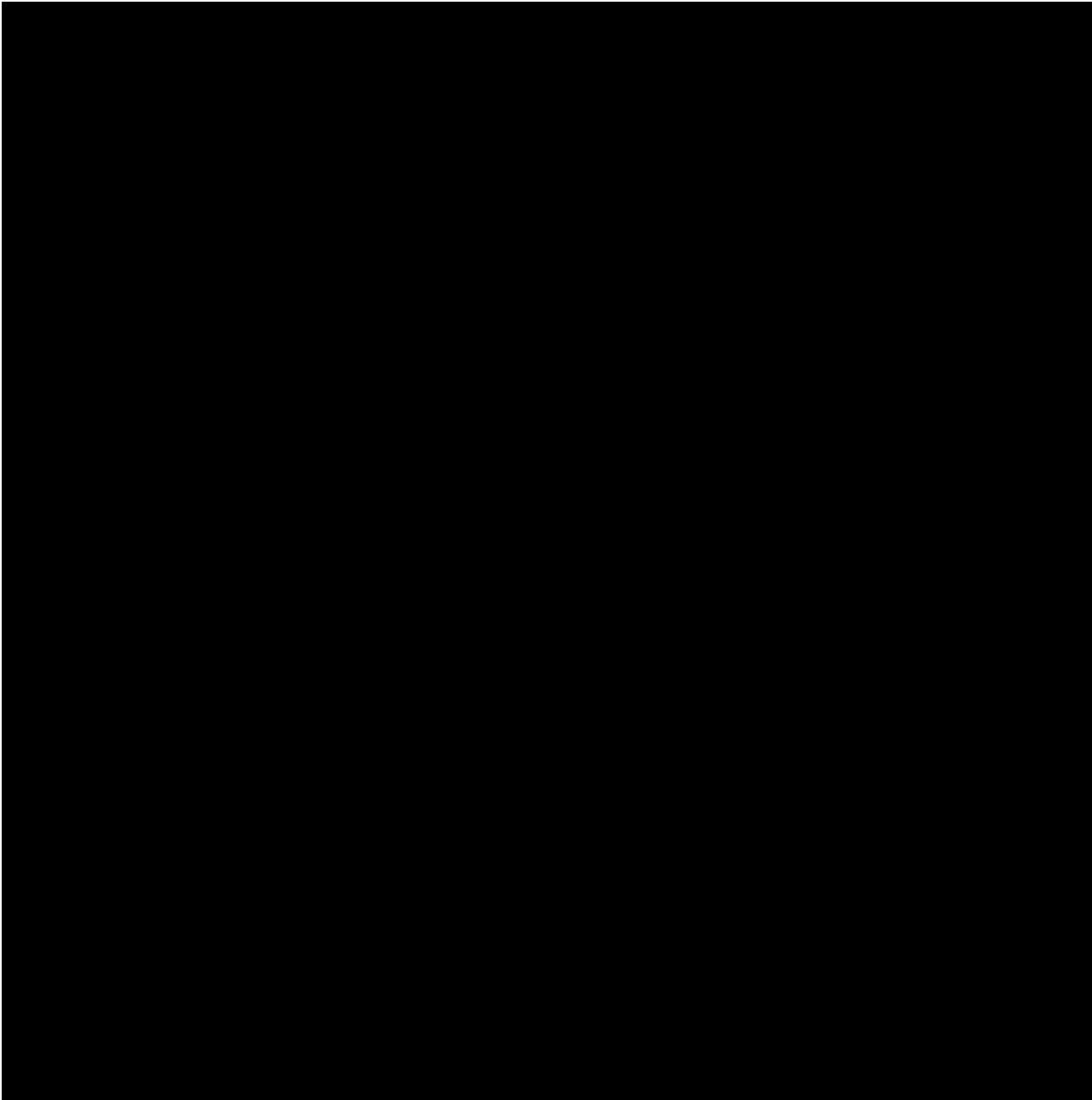
2 121. The fourth example also demonstrates Copilot suggesting multiple modified
3 copies of code written by Doe 5 in response to a sequence of prompts, which is a common way of
4 using Copilot. To protect Doe 5's identity, the paragraphs describing the code will be redacted.

5 122. [REDACTED]

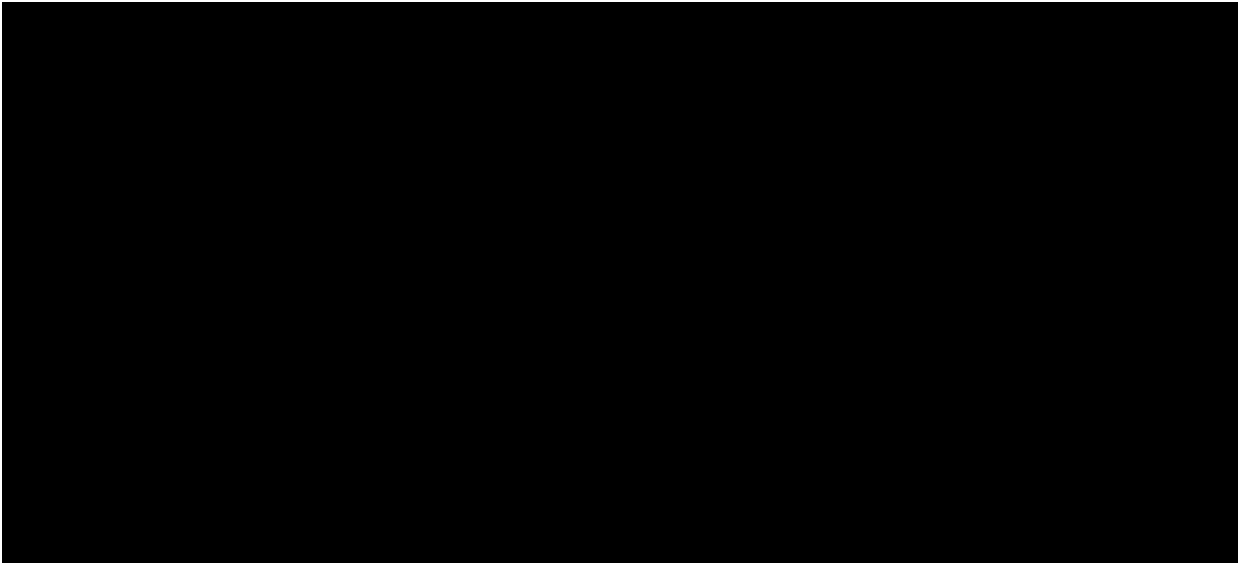
6 [REDACTED] subject to the MIT License. [REDACTED]

7 [REDACTED]. The first three tests from the original source file are shown

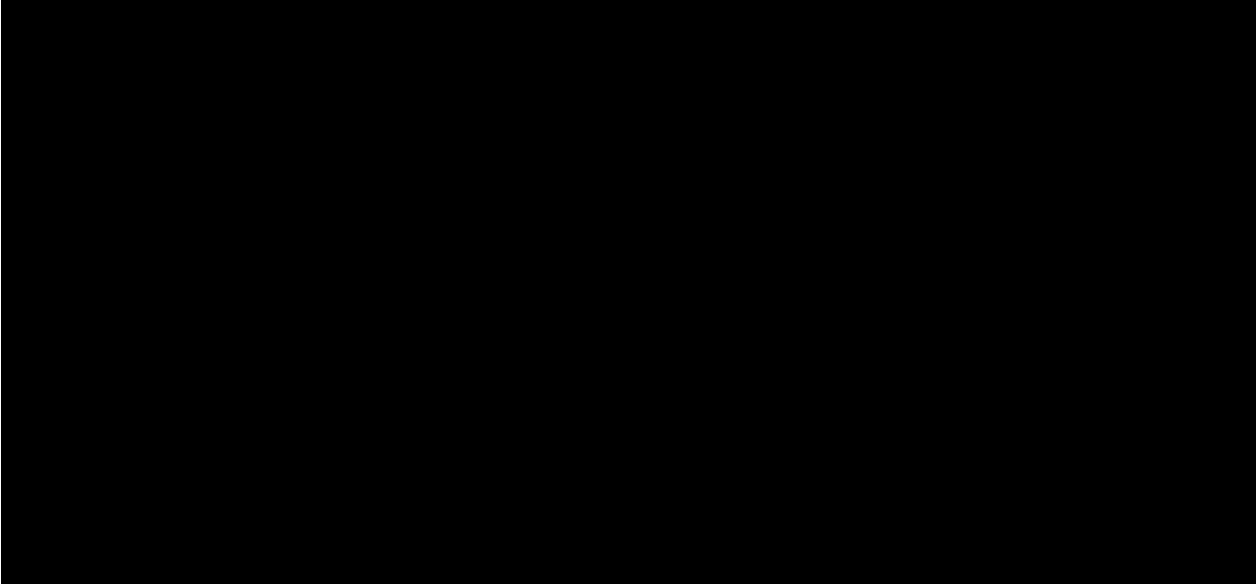
8 below:



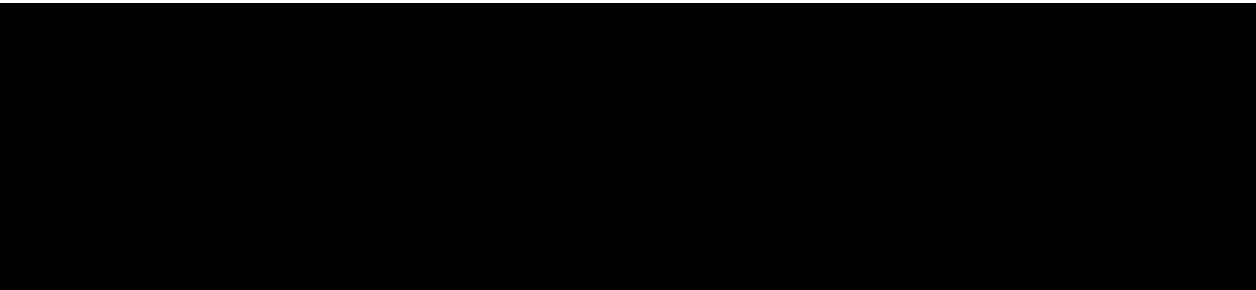
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28



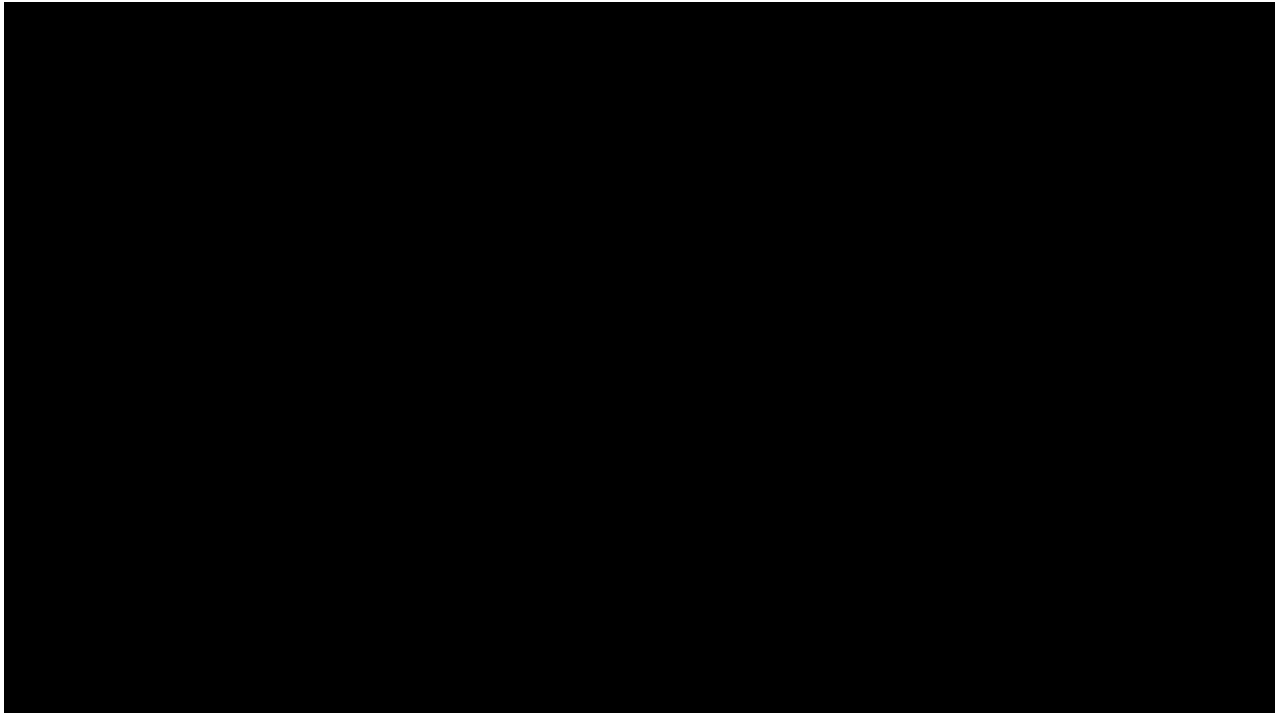
123. When Copilot is prompted with the first section of Doe 5's code, comprising the first complete test and the name of the second:



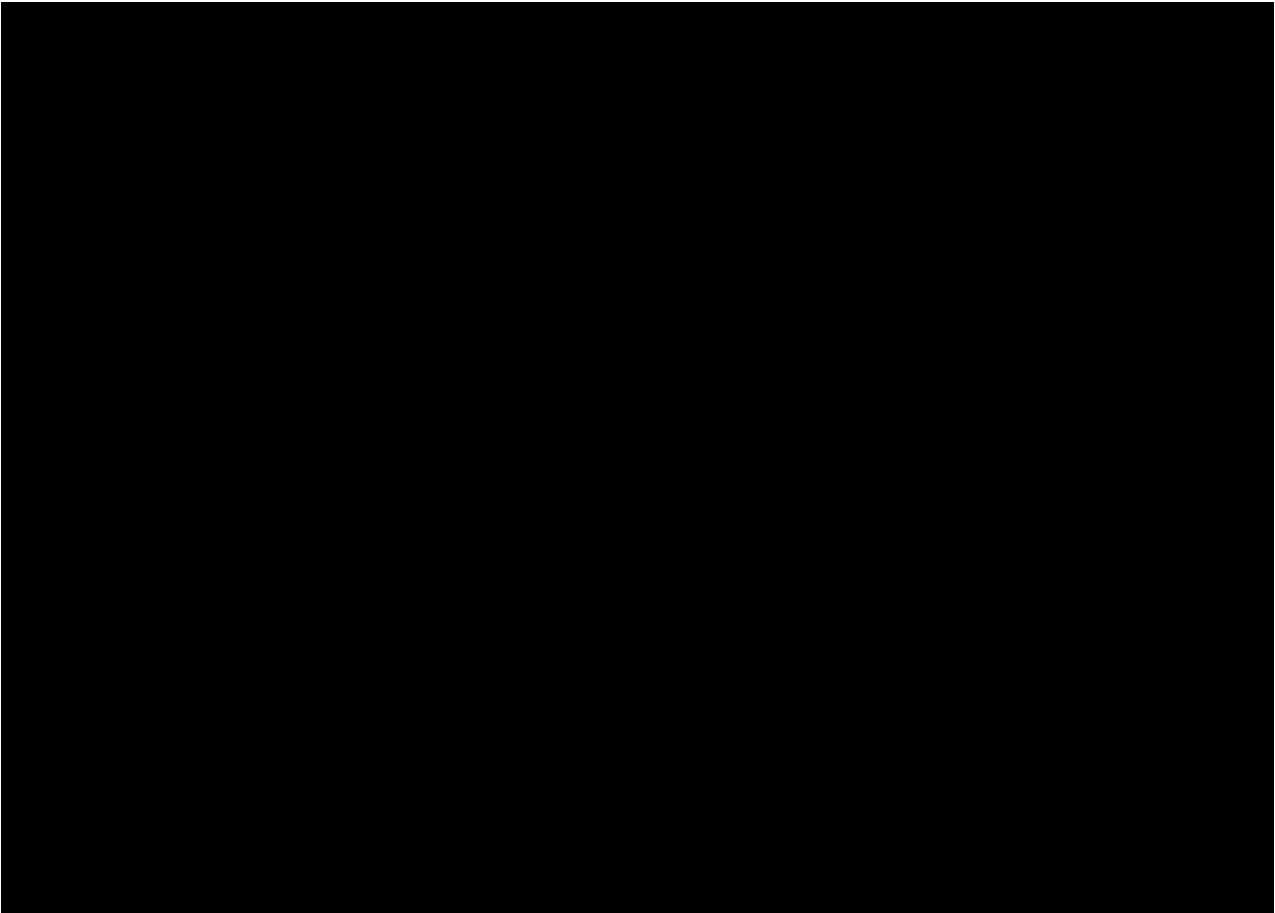
The first suggestion from Copilot offers to complete the second test with a verbatim copy of Doe 5's original code:



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



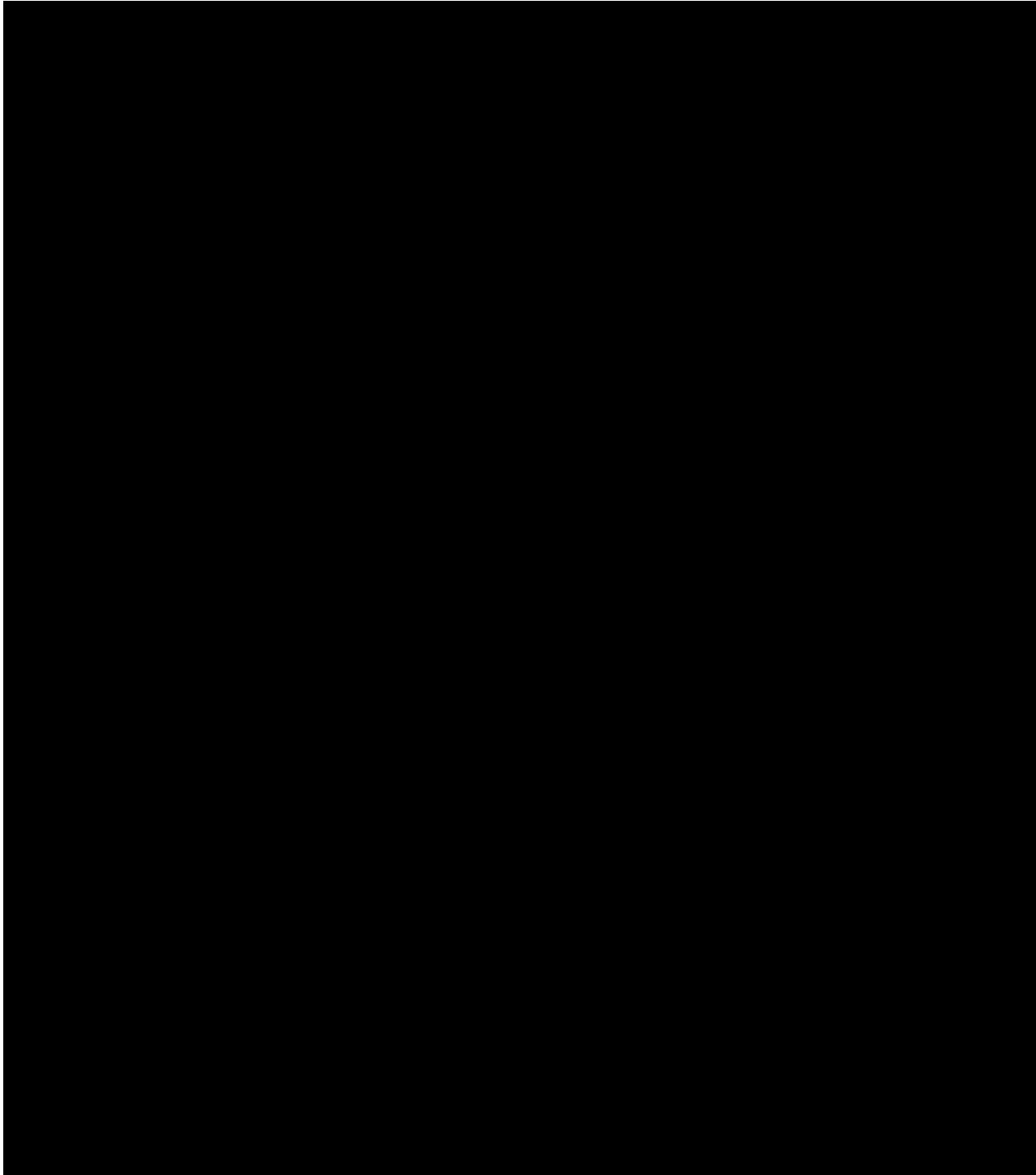
124. When Copilot’s suggestion is accepted and the name of Doe 5’s third test is appended, the next prompt to Copilot looks like this:

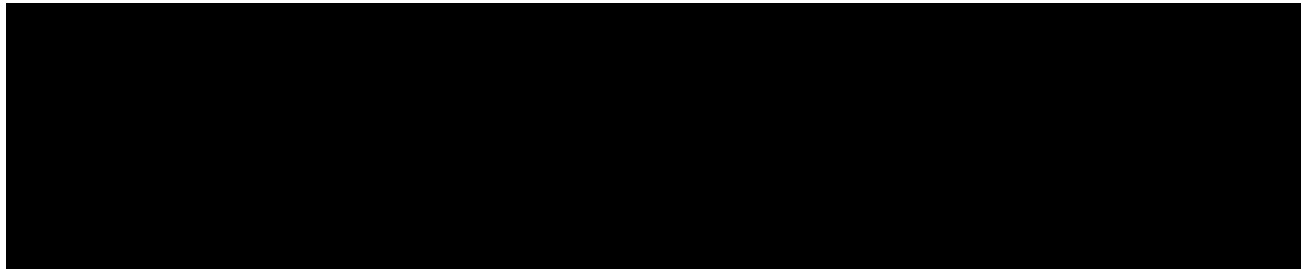


1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28



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





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

126. Because Copilot is (repeatedly) reproducing Doe 5’s code essentially verbatim, the Copilot suggestions need to follow the requirements of Doe 5’s license (the MIT License) for that code, including providing attribution. They do not. Copilot also did not reproduce Doe 5’s license.

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

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

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

129. Codex and Copilot have no way to determine whether license text or other Copyright Management Information (“CMI”)²¹ is part of the code it appears immediately before

²¹ CMI is defined in detail below in Paragraph 187.

1 or after. Unless instructed otherwise, it will assume that CMI that usually appears just before a
2 given block of code is an important part of that code or otherwise necessary for it to function.

3 130. It is a common practice to provide the applicable license text at the top of every
4 source file in the codebase. The purpose of this practice is to avoid the code from being divorced
5 from the license. This may occur via “vendoring,” a method of creating a derivative work by
6 including source files from a copyrighted project directly into another project without following
7 the terms of the license or providing attribution or a copyright notice. Copilot circumvents this
8 protective measure to mask the degree of vendoring it engages in.

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

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

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

27
28 ²² <https://github.blog/2021-06-30-github-copilot-research-recitation/>.

1 134. In GitHub’s hands, the propensity for small cosmetic variations in Copilot’s
2 Output is a feature, not a bug. These small cosmetic variations mean that GitHub can deliver to
3 Copilot customers unlimited modified copies of Licensed Materials without ever triggering
4 Copilot’s verbatim-code filter. AI models like Copilot often have a setting called *temperature* that
5 specifically controls the propensity for variation in their output. On information and belief,
6 GitHub has optimized the temperature setting of Copilot to produce small cosmetic variations of
7 the Licensed Materials as often as possible, so that GitHub can deliver code to Copilot users that
8 *works* the same way as verbatim code, while claiming that Copilot only produces verbatim code 1%
9 of the time. Copilot is an ingenious method of software piracy.

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

13 GitHub will defend you against any claim by an unaffiliated third-
14 party that your use of GitHub Copilot misappropriated a trade
15 secret or directly infringes a patent, copyright, trademark, or other
16 intellectual property right of a third party, up to the greater of
17 \$500,000.00 USD or the total amount paid to GitHub for the use
18 of GitHub Copilot during the 12 months preceding the claim.
19 GitHub’s defense obligations do not apply if (i) the claim is based
20 on Code that differs from a Suggestion provided by GitHub
21 Copilot, (ii) you fail to follow reasonable software development
22 review practices designed to prevent the intentional or inadvertent
23 use of Code in a way that may violate the intellectual property or
24 other rights of a third party, or (iii) you have not enabled all filtering
25 features available in GitHub Copilot.

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

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

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

1 138. 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.

4 139. 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.

8 140. 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.

12 141. 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.²³

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.

22 142. 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 ²³ <https://www.gnu.org/licenses/old-licenses/gpl-1.0.en.html>.

1 143. 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.

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

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

7 145. 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.

13 146. 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.

15 147. 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.

18 148. 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.

22 149. 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.

24 150. 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.”²⁴

3 151. 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.”²⁵
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.

9 152. In 2001, Microsoft was the defendant in a major software-related antitrust case,
10 *United States v. Microsoft Corporation*.²⁶ 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.”²⁷

17 153. 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.²⁸

20 154. 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 ²⁴ <http://www.catb.org/esr/halloween/halloween1.html>.

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

26 ²⁶ No. Civ.A. 00-1457 TPJ.

27 ²⁷ *Jackson v. Microsoft Corp.*, 135 F. Supp. 2d 38 (D.D.C. 2001).

28 ²⁸ <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 **J. GitHub Was Designed to Cater to Open-Source Projects**

6 155. 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.

9 156. 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.

12 157. 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.

15 158. 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.

19 159. 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.

24 160. 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.

1 161. 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.

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

12 163. 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.

17 164. 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.

21 165. 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.

23 166. 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 **K. OpenAI Is Intertwined with Microsoft and GitHub**

4 167. 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.

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

15 169. 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.”

17 170. 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.

19 171. 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.”²⁹

24 172. 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 ²⁹ [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/).

1 173. 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 176). 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.”

8 174. 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.³⁰ 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.

13 175. 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.”³¹

17 176. 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.

21 177. 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.

24 178. 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 ³⁰ <https://www.microsoft.com/en-us/Investor/earnings/FY-2023-Q1/press-release-webcast/>.

28 ³¹ <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.

3 179. In another joint project, Microsoft and OpenAI recently launched a preview of a
4 product called “Azure OpenAI Service.”³² 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 **L. Conclusion of Factual Allegations**

11 180. 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.

15 181. 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.”³³

21 182. 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 ³² <https://azure.microsoft.com/en-us/products/cognitive-services/openai-service/>.

28 ³³ https://www.digibarn.com/collections/newsletters/homebrew/V2_01/gatesletter.html

VIII. CLAIMS FOR RELIEF

**COUNT 1
VIOLATION OF THE DIGITAL MILLENNIUM COPYRIGHT ACT
17 U.S.C. §§ 1201–1205
(Direct, Vicarious, and Contributory)
(Against All Defendants)**

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

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

185. As described herein, Defendants have distributed copies of Plaintiffs’ code knowing that CMI has been removed or altered while knowing or having reasonable grounds to know that it will induce, enable, facilitate, or conceal infringement in violation of Section 1202(b)(3) of the DMCA.

186. 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.

187. 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

1 f. identifying numbers or symbols referring to CMI or links to CMI.

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

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

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

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

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

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

21 194. Defendants had access to but were not licensed by Plaintiffs nor the Class to create
22 Derivative Works³⁴ based upon the Licensed Materials.

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

25
26
27 ³⁴ "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.

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

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

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

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

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

26 201. By omitting, altering and/or concealing CMI from Copilot's Output, Defendants
27 have reasonable grounds to know that innocent infringers are induced or enabled to copy the
28 Licensed Materials, because CMI has been removed. Without the CMI, Defendants have

1 reasonable grounds to know copyright infringement is facilitated or concealed, because Plaintiffs
2 and the Class have the difficult or impossible task of proving the Licensed Materials belong to
3 them.

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

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

13 204. Defendants did not request or obtain permission from Plaintiffs and the Class to
14 use the Licensed Materials for Defendants' Copilot product.

15 205. Defendants use of the Licensed Materials does not follow the requirements of the
16 Suggested Licenses associated with the Licensed Materials. In particular, Copilot fails to provide
17 attribution for the creator nor the owner of the Work. Copilot fails to include the required
18 copyright notice included in the License. Copilot fails to include the applicable Suggested
19 License's text.

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

23 **We protect our intellectual property investments in a variety of**
24 **ways. We work actively in the U.S. and internationally to**
25 **ensure the enforcement of copyright, trademark, trade secret,**
26 **and other protections that apply to our software and hardware**
27 **products, services, business plans, and branding.** We are a
28 leader among technology companies in pursuing patents and
currently have a portfolio of over 69,000 U.S. and international
patents issued and over 19,000 pending worldwide. While we
employ much of our internally-developed intellectual property
exclusively in our products and services, we also engage in

1 outbound licensing of specific patented technologies that are
2 incorporated into licensees' products. From time to time, we enter
3 into broader cross-license agreements with other technology
4 companies covering entire groups of patents. We may also purchase
5 or license technology that we incorporate into our products and
6 services. At times, we make select intellectual property broadly
7 available at no or low cost to achieve a strategic objective, such as
8 promoting industry standards, advancing interoperability,
9 supporting societal and/or environmental efforts, or attracting and
10 enabling our external development community. **Our increasing
11 engagement with open source software will also cause us to
12 license our intellectual property rights broadly in certain
13 situations.**

14 Microsoft Corporation Annual Report, Form 10-K at 27 (July 28, 2022) (emphasis added).³⁵

15 207. GitHub, which offers the Copilot product jointly with OpenAI, also has extensive
16 experience with the DMCA. GitHub knows or reasonably should know that the Licensed
17 Materials it hosts are subject to copyright. It provides the language of the Suggested Licenses to
18 users, all of which include copyright notices. Its 2022 Transparency Report—January to June³⁶
19 states: “Copyright-related takedowns (which we often refer to as DMCA takedowns) are
20 particularly relevant to GitHub because so much of our users’ content is software code and can be
21 eligible for copyright protection.”³⁷ In the first six months of 2022, GitHub processed 1220
22 DMCA takedown requests. Its DMCA Takedown Policy³⁸ notes “GitHub probably never would
23 have existed without the DMCA.”

24 208. GitHub also knows or reasonably should know the portions of the DMCA giving
25 rise to Plaintiffs’ claim. In its 2021 Transparency Report, “Before removing content based on
26 alleged circumvention of copyright controls (under Section 1201 of the US DMCA or similar laws
27 in other countries), we carefully review both the legal and technical claims, and we sponsor a
28 Developer Defense Fund to provide developers with meaningful access to legal resources.”³⁹

24 ³⁵ <https://microsoft.gcs-web.com/static-files/07cf3c30-cfc3-4567-b20f-f4b0f0bd5087/>.

25 ³⁶ <https://github.blog/2022-08-16-2022-transparency-report-january-to-june/>.

26 ³⁷ <https://github.blog/2022-08-16-2022-transparency-report-january-to-june/>.

27 ³⁸ <https://docs.github.com/en/site-policy/content-removal-policies/dmca-takedown-policy#what-is-the-dmca/>.

28 ³⁹ <https://github.blog/2022-01-27-2021-transparency-report/>.

1 209. GitHub is aware that Copilot’s removal of CMI is illegal. For example, it states
2 that “publishing or sharing tools that enable circumvention are not [permitted]”⁴⁰ and
3 “Distributing tools that enable circumvention is prohibited, even if their use by developers falls
4 under the exemption [for security research].”⁴¹ GitHub has also frequently published articles
5 discussing the DMCA, its application, and the Copyright Office’s guidance on its scope and
6 exceptions.⁴²

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

15 211. Plaintiffs and the Class are further entitled to recover from Defendants the actual
16 or statutory damages Plaintiffs and the Class sustained pursuant to 17 U.S.C. § 1203(c) and for
17 Plaintiffs’ and the Class’s costs and attorneys’ fees in enforcing the Licenses. Plaintiffs and the
18 Class are also entitled to recover as restitution from Defendants for any unjust enrichment,
19 including gains, profits, and advantages that Defendants have obtained as a result of their breach
20 of the Licenses.

21 212. Defendants conspired together and acted jointly and in concert pursuant to their
22 scheme to commit the acts that violated the DMCA alleged herein.

23 213. Defendants induced Copilot users to unknowingly violate the DMCA by
24 withholding attribution, licensing, and other information as described herein.

25
26 ⁴⁰ <https://github.blog/2020-11-19-take-action-dmca-anti-circumvention-and-developer-innovation/#what-dmca-exemptions-do-not-do/>.

27 ⁴¹ <https://github.blog/2021-11-23-copyright-office-expands-security-research-rights/>.

28 ⁴² *See, e.g.*, Footnotes 36–41.

COUNT 2
BREACH OF CONTRACT—OPEN-SOURCE LICENSE VIOLATIONS
California Common Law
(Against All Defendants)

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

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

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

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

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

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

220. Plaintiffs and the Class have appended one of the Suggested Licenses to each of the Licensed Materials.

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

1 222. Defendants have substantially and materially breached the applicable Licenses by
2 failing to provide the source code of Copilot nor a written offer to provide the source code upon
3 the request of each licensee.

4 223. Defendants have substantially and materially breached the applicable Licenses by
5 failing to provide attribution to the creator and/or owner of the Licensed Materials.

6 224. Defendants have substantially and materially breached the applicable Licenses by
7 failing to include copyright notices when Copilot Outputs copyrighted OS code.

8 225. Defendants have substantially and materially breached the applicable Licenses by
9 failing to identify the License applicable to the Work and/or including its text when Copilot
10 Outputs code including a portion of a Work.

11 226. Plaintiffs and the Class have suffered monetary damages as a result of Defendants'
12 conduct.

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

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

20 229. Plaintiffs and the Class are further entitled to recover from Defendants the
21 damages Plaintiffs and the Class sustained—including consequential damages—for Plaintiffs' and
22 the Class's costs in enforcing their contractual rights. Plaintiffs and the Class are also entitled to
23 recover as restitution from Defendants for any unjust enrichment, including gains, profits, and
24 advantages that Defendants have obtained as a result of their breach of contract.

COUNT 3
BREACH OF CONTRACT — SELLING LICENSED MATERIALS
IN VIOLATION OF GITHUB’S POLICIES
California Common Law
(Against GitHub)

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

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

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

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

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

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

236. GitHub has also breached the implied covenant of good faith and fair dealing. GitHub has long held itself out as a good citizen of the global open-source community. GitHub’s Policies were designed to attract Plaintiffs and the Class to become users of the GitHub website

1 by supporting their open-source efforts with fair and ethical practices. By releasing Copilot,
2 GitHub created a product designed to compete with Plaintiffs and the Class and undermine their
3 individual open-source communities. In so doing, GitHub did not act fairly or in good faith.

4 237. Plaintiffs and the Class have suffered monetary damages as a result of GitHub's
5 conduct.

6 238. GitHub's conduct is causing and, unless enjoined and restrained by this Court,
7 will continue to cause Plaintiffs and the Class great and irreparable injury that cannot fully be
8 compensated or measured in money.

9 239. As a direct and proximate result of these material breaches by GitHub, Plaintiffs
10 and the Class are entitled to an injunction requiring GitHub to comply with all the terms of the
11 GitHub Policies.

12 240. Plaintiffs and the Class are further entitled to recover from GitHub the damages
13 Plaintiffs and the Class sustained—including consequential damages—for Plaintiffs' and the
14 Class's costs in enforcing GitHub's Policies. Plaintiffs and the Class are also entitled to recover as
15 restitution from GitHub for any unjust enrichment, including gains, profits, and advantages that it
16 has obtained as a result of its breaches of the GitHub Policies.

17
18 **COUNT 4**
19 **INTENTIONAL INTERFERENCE**
20 **WITH PROSPECTIVE ECONOMIC RELATIONS**
21 **California Common Law**
22 **(Against All Defendants)**

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

25 242. Open-source software programmers invite other programmers to view, use and
26 modify their code and to make changes and improvements to it subject to requirements in certain
27 licenses. When a programmer uses an open-source software, a contract is formed based on the
28 terms of the particular open-source license.

29 243. Although no money changes hands in open-source licensing, courts have long
30 recognized that there are substantial benefits, including economic benefits, to the creation and

1 distribution of open-source code subject to these open-source licenses. For example, program
2 creators can generate market share for their programs, increase their reputation both nationally
3 and internationally, and discover new improvements to their open-source project. These benefits,
4 however, rely on the proliferation of the licenses and obligations that come with the license along
5 with the code that is subject to that license.

6 244. GitHub was founded by open-source programmers and has long held itself out as a
7 gathering point for the global open-source community. For open-source programmers, including
8 Plaintiffs and the Class, part of the benefit of becoming a GitHub customer and sharing code
9 there was to make it easier for members of the global open-source community to discover their
10 work, and thereby accrue a user community of contributors and collaborators specific to their
11 open-source projects.

12 245. User communities create the probability of future economic benefit in a number of
13 ways. Users can provide bug reports, saving authors from having to discover every bug
14 themselves. Users can provide new code that fixes bugs or adds features, saving authors from
15 having to write every line of code themselves. Users sometimes arrange financial contracts with
16 authors for extra licensing rights, or custom features, or technical support. The exposure from a
17 user community can also bring collateral benefits, like job offers or research grants.

18 246. Plaintiffs and Class members chose to become GitHub customers specifically to
19 avail themselves of these benefits of GitHub and optimize the likelihood of accruing communities
20 of other GitHub customers for their own projects, and to benefit from the future economic
21 benefits likely to arise from those relationships.

22 247. In other words, Plaintiffs and Class members posted their code on GitHub with an
23 expectation that other programmers would use, modify, copy or otherwise iterate on their posted
24 code subject to the terms of the open-source licenses the code was published subject to.

25 248. GitHub's status as a focal point of the global open-source community was one of
26 the main reasons Microsoft wanted to own it. At the time it acquired GitHub, Microsoft pledged
27 to uphold these virtues. But Defendants' project of harvesting mass quantities of public open-
28 source code on GitHub for training Codex and Copilot represented an inversion of these

1 priorities. Codex and Copilot essentially act as walled gardens that provide an alternative
2 interface to the same open-source code, a process sometimes called “vendorization”.

3 249. By failing to provide information about the Suggested Licenses attached to the
4 Licensed Materials, Defendants intentionally prevented Copilot users from becoming part of the
5 user communities that would ordinarily accrete around the open-source projects of Plaintiffs and
6 the Class. Instead, Defendants reserved those benefits for themselves.

7 250. Defendants knew that they were interfering with Plaintiffs and Class members’
8 prospective open-source relationships because Defendants knew that Codex and Copilot were
9 emitting code subject to open-source licenses without the licenses attached.

10 251. Defendants have therefore intentionally and wrongfully interfered with the
11 prospective business interests and expectations of Plaintiffs and the Class.

12 252. Plaintiffs have been deprived of the economic benefits of open-source licenses.
13 Plaintiffs and the Class have suffered monetary, reputational, and other damages as a result of
14 Defendants’ conduct.

15 253. Unless enjoined and restrained by this Court, Defendants’ conduct will continue
16 to cause Plaintiffs and the Class great and irreparable injury that cannot fully be compensated or
17 measured in money.

18
19 **COUNT 5**
20 **NEGLIGENT INTERFERENCE**
21 **WITH PROSPECTIVE ECONOMIC RELATIONS**
California Common Law
(Against All Defendants)

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

24 255. As described in Paragraphs 242 through 247 herein, Plaintiffs and Class members
25 post their code subject to open-source licenses in order to enjoy the economic benefits associated
26 with the distribution of open-source software code.

27 256. Defendants knew or should have known that Plaintiffs and the Class chose to
28 become GitHub customers specifically to optimize the likelihood of accruing communities of

1 other GitHub customers for their own projects, and to benefit from the future economic benefits
2 likely to arise from those open-source licensing relationships.

3 257. Defendants knew or should have known that by scraping the Licensed Materials of
4 Plaintiffs and the Class and using it to create competing products (namely, Codex and Copilot)
5 that did not adhere to the obligations of the Suggested Licenses, that Defendants would disrupt
6 the formation and growth of these user communities, and also disrupt the economic benefits that
7 would ordinarily accrue to Plaintiffs and the Class from the growth of those user communities.

8 258. Defendants' conduct falls far outside the boundaries of fair competition because
9 Defendants have leveraged GitHub's status as a gathering point for the global open-source
10 community to undermine that very community, including Plaintiffs and the Class. Defendants
11 have leveraged the permissive nature of the Suggested Licenses to undermine the Licensed
12 Materials, thereby harming Plaintiffs and the Class.

13 259. Defendants' conduct could not have been performed without a direct effect on
14 Plaintiffs' economic interests as Plaintiffs' open-source code was hosted on GitHub and used to
15 train Codex and Copilot.

16 260. Plaintiffs uploaded their open-source software code on GitHub with the
17 reasonable expectation that other open-source programmers would use, modify, copy or
18 otherwise iterate on their code.

19 261. The adverse effect on Plaintiffs was foreseeable as Defendants are aware of the
20 obligations that carry with open-source licenses. Further, Defendants knew that Codex and
21 Copilot copied code used for training the models without the associated licenses attached and
22 without any of the necessary obligations that carry with the licenses, e.g., attribution.

23 262. Even after Defendants knew that Codex and Copilot would reproduce code
24 verbatim, including in some instances, with the incorrect licenses or license text, Defendants
25 continued to operate Codex and Copilot. In other words, Defendants continued to operate Codex
26 and Copilot after they should have known that Codex and Copilot's operation were depriving
27 Plaintiffs of the economic benefits of open-source distribution.

1 emitted without the proper licensing and attribution required by the licenses Plaintiffs and Class
2 members publish their material subject to.

3 280. The unlawful business practices described herein violate the common law because
4 Plaintiffs and Class members invested substantial time, skill and money in developing the
5 Licensed Materials. Defendants misappropriated and used Plaintiffs and Class members'
6 Licensed materials without authorization or consent in order to, *inter alia*, train and develop
7 Codex and Copilot.

8 281. Plaintiffs and the Class have suffered economic injury as a result of Defendants'
9 conduct. Specifically, there are economic benefits to the creation of open-source works such as
10 generating market share for programs, increasing national or international reputation by
11 incubating open-source projects, and deriving value from improvements to software based on
12 suggestions by end-users.

13 **COUNT 8**
14 **NEGLIGENCE**
15 **California Common Law**
(Against All Defendants)

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

18 283. Defendants owed a duty of reasonable care toward Plaintiffs and the Class based
19 upon Defendants' relationship to them. This duty is based upon Defendants' contractual
20 obligations, custom and practice, right to control information in its possession, exercise of control
21 over the information in its possession, authority to control the information in its possession, and
22 the commission of affirmative acts that resulted in said harms and losses. Additionally, this duty is
23 based on the requirements of California Civil Code section 1714 requiring all "persons,"
24 including Defendants, to act in a reasonable manner toward others.

25 284. Defendants breached their duties by negligently, carelessly, and recklessly
26 collecting, maintaining, and controlling their customers' Licensed Materials and engineering,
27 designing, maintaining, and controlling systems—including Codex and Copilot—which are
28 trained on Plaintiffs' and Class members' Licensed Materials without their authorization.

1 285. Microsoft and GitHub owed its users a duty of care to safeguard and maintain
2 Licensed Materials on its website and to prevent unauthorized use of the Licensed Materials.

3 286. Microsoft and GitHub also owed its user a duty of care not to itself use the
4 Licensed Materials in a way that would foreseeably cause Plaintiffs and Class members injury, for
5 instance, by using Licensed Materials to train Copilot.

6 287. OpenAI owed Plaintiffs and Class members a duty of care by using open-source
7 code in violation of open-source licenses to train Codex and Copilot.

8 288. Defendants, through their unlawful acts described herein, breach of their duties
9 proximately caused Plaintiffs and Class members injuries.

10 289. Plaintiffs and Class members' injuries were foreseeable because Defendants are
11 aware of the benefits of open-source licensing because they hold themselves out as advocates for
12 open source, and serve as platforms to host open-source software.

13 **IX. DEMAND FOR JUDGMENT**

14 **WHEREFORE**, Plaintiffs requests that the Court enter judgment on their behalf and
15 on behalf of the Class defined herein, by adjudging and decreeing that:

16 290. This action may proceed as a class action, with Plaintiffs serving as Class
17 Representatives, and with Plaintiffs' counsel as Class Counsel;

- 18 a) Judgment in favor of Plaintiffs and the Class and against Defendants;
- 19 b) Permanent injunctive relief, including but not limited to making changes to its
20 Copilot product to ensure that all applicable information set forth in 17 U.S.C. §
21 1203(b)(1) is included in along with any Output including associated code;
- 22 c) An order of costs and allowable attorney's fees pursuant to 17 U.S.C. §
23 1203(b)(4)–(5);
- 24 d) An award of statutory damages pursuant to 17 U.S.C. § 1203(b)(3) and 17 U.S.C. §
- 25
- 26
- 27
- 28

1 1203(c)(3),⁴³ or, in the alternative, an award of actual damages and any additional
2 profits pursuant to 17 U.S.C. § 1203(c)(2) (including tripling damages pursuant to
3 17 U.S.C. § 1203(c)(4) if applicable);

4 e) An award of damages for harms resulting from Defendants' breach of Licenses;

5 f) An award of damages, including punitive damages, for harms resulting from
6 Defendants' tortious interference in Plaintiffs' and the Class's prospective
7 economic advantage;

8 g) An award of damages in the amount Defendants have been unjustly enriched
9 through their conduct as alleged herein as well as punitive damages in connection
10 with this conduct;

11 h) An award of damages, including punitive damages, for harms resulting from
12 Defendants' acts of unfair competition;

13 i) An award of damages for harms resulting from GitHub's breach of the GitHub
14 Policies; and

15 j) An award of damages, including punitive damages, for harms resulting from
16 Defendants' negligence including in the failure to control Plaintiffs' and the
17 Class's Licensed Materials.

18 291. Injunctive relief sufficient to alleviate and stop Defendants' unlawful conduct
19 alleged herein.

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

23 ⁴³ 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
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