

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

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

14 Laura M. Matson (*pro hac vice* pending)  
15 **LOCKRIDGE GRINDAL NAUEN PLLP**  
16 100 Washington Avenue South, Suite 2200  
17 Minneapolis, MN 55401  
18 Telephone: (612) 339-6900  
19 Facsimile: (612) 339-0981  
20 Email: lmmatson@locklaw.com

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

23 **UNITED STATES DISTRICT COURT**  
24 **NORTHERN DISTRICT OF CALIFORNIA**  
25 **SAN FRANCISCO DIVISION**

26 **Stewart O’Nan**, an individual;  
27 **Abdi Nazemian**, an individual; and  
28 **Brian Keene**, an individual;

Individual and Representative Plaintiffs,

v.

**Databricks, Inc.**, a Delaware corporation; and  
**Mosaic ML, Inc.**, a Delaware corporation;

Defendants.

Case No.

**COMPLAINT**

**CLASS ACTION**

**DEMAND FOR JURY TRIAL**

1 Plaintiffs Stewart O’Nan, Abdi Nazemian, and Brian Keene (together “Plaintiffs”), on behalf of  
2 themselves and all others similarly situated, bring this class-action complaint (“Complaint”) against  
3 defendants Mosaic ML, Inc. (“MosaicML”) and Databricks, Inc. (“Databricks”) (together  
4 “Defendants”).

## 6 OVERVIEW

7 1. *Artificial intelligence*—commonly abbreviated “AI”—denotes software that is designed  
8 to algorithmically simulate human reasoning or inference, often using statistical methods.

9 2. A *large language model* is an AI software program designed to emit convincingly  
10 naturalistic text outputs in response to user prompts. MosaicML Pretrained Transformer (“MPT”) is  
11 a series of large language models created by MosaicML and distributed by Databricks.

12 3. Rather than being programmed in the traditional way—that is, by human programmers  
13 writing code—a large language model is *trained* by copying an enormous quantity of textual works,  
14 extracting protected expression from these works, and transforming that protected expression into a  
15 large set of numbers called *weights* that are stored within the model. These weights are entirely and  
16 uniquely derived from the protected expression in the training dataset. Whenever a large language  
17 model generates text output in response to a user prompt, it is performing a computation that relies on  
18 these stored weights, with the goal of imitating the protected expression ingested from the training  
19 dataset.

20 4. Plaintiffs and Class members are authors. They own registered copyrights in certain  
21 books that were included in the training dataset that MosaicML has admitted copying to train its MPT  
22 models. Plaintiffs and Class members never authorized MosaicML to use their copyrighted works as  
23 training material.

24 5. MosaicML copied these copyrighted works multiple times to train its MPT models.

25 6. Databricks, as the corporate parent of MosaicML and distributor of the MPT models,  
26 has also commercially benefitted from these acts of massive copyright infringement.

**JURISDICTION AND VENUE**

7. This Court has subject-matter jurisdiction under 28 U.S.C. § 1331 because this case arises under the Copyright Act (17 U.S.C. § 501).

8. Jurisdiction and venue are proper in this judicial district under 28 U.S.C. § 1391(c)(2) because Defendants are headquartered in this district. Defendants created the MPT language models and distribute them commercially. Therefore, a substantial part of the events giving rise to the claim occurred in this District. A substantial portion of the affected interstate trade and commerce was carried out in this District. Defendants have transacted business, maintained substantial contacts, and/or committed overt acts in furtherance of the illegal scheme and conspiracy throughout the United States, including in this District. Defendants’ conduct has had the intended and foreseeable effect of causing injury to persons residing in, located in, or doing business throughout the United States, including in this District.

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

**PLAINTIFFS**

10. Plaintiff Stewart O’Nan is an author who lives in Pennsylvania. Mr. O’Nan owns registered copyrights in multiple books, including *Last Night at the Lobster*.

11. Plaintiff Abdi Nazemian is an author who lives in California. Mr. Nazemian owns registered copyrights in multiple books, including *Like a Love Story*.

12. Plaintiff Brian Keene is an author who lives in Pennsylvania. Mr. Keene owns registered copyrights in multiple books, including *Ghost Walk*.

13. A non-exhaustive list of registered copyrights owned by Plaintiffs is included as Exhibit A.

1 **DEFENDANTS**

2 14. Defendant Databricks is a Delaware corporation with its principal place of business at  
3 160 Spear Street, 13th Floor, San Francisco CA 94105. Databricks acquired MosaicML in July 2023.

4 15. Defendant MosaicML is a Delaware corporation with its principal place of business at  
5 501 2nd Street, Suite 202, San Francisco CA 94107. MosaicML operates as a subsidiary of Databricks.

6  
7 **AGENTS AND CO-CONSPIRATORS**

8 16. The unlawful acts alleged against Defendants in this class action complaint were  
9 authorized, ordered, or performed by the Defendants’ respective officers, agents, employees,  
10 representatives, or shareholders while actively engaged in the management, direction, or control of the  
11 Defendants’ businesses or affairs. The Defendants’ agents operated under the explicit and apparent  
12 authority of their principals. Each Defendant, and its subsidiaries, affiliates, and agents operated as a  
13 single unified entity.

14 17. Various persons or firms not named as defendants may have participated as co-  
15 conspirators in the violations alleged herein and may have performed acts and made statements in  
16 furtherance thereof. Each acted as the principal, agent, or joint venture of Defendants with respect to  
17 the acts, violations, and common course of conduct alleged herein.

18  
19 **FACTUAL ALLEGATIONS**

20 18. MosaicML was founded in 2020 to provide tools to facilitate the training of AI models.

21 19. In May 2023, MosaicML released the first in its MPT series of *large language models*,  
22 called MPT-7B. A large language model (“LLM”) is AI software designed to emit convincingly  
23 naturalistic text outputs in response to user prompts.

24 20. Though an LLM is a software program, it is not created the way most software  
25 programs are—that is, by human software programmers writing code. Rather, an LLM is *trained* by  
26 copying an enormous quantity of textual works and then feeding these copies into the model. This  
27 corpus of input material is called the *training dataset*.

1           21.     During training, the LLM copies and ingests each textual work in the training dataset  
2 and extracts protected expression from it. The LLM progressively adjusts its output to more closely  
3 approximate the protected expression copied from the training dataset. The LLM records the results of  
4 this process in a large set of numbers called *weights* that are stored within the model. These weights are  
5 entirely and uniquely derived from the protected expression in the training dataset. For instance, the  
6 MPT-7B language model is so named because the model stores seven billion (“7B”) weights derived  
7 from protected expression in its training dataset.

8           22.     Once the LLM has copied and ingested the textual works in the training dataset and  
9 transformed the protected expression into stored weights, the LLM is able to emit convincing  
10 simulations of natural written language in response to user prompts. Whenever an LLM generates text  
11 output in response to a user prompt, it is performing a computation that relies on these stored weights,  
12 with the goal of imitating the protected expression ingested from the training dataset.

13           23.     Much of the material in MosaicML’s training dataset, however, comes from copyrighted  
14 works—including books written by Plaintiffs and Class members—that were copied by MosaicML  
15 without consent, without credit, and without compensation.

16           24.     In May 2023, MosaicML first announced MPT-7B in a blog post called “Introducing  
17 MPT-7B: A New Standard for Open-Source, Commercially Usable LLMs.”<sup>1</sup> In the blog post,  
18 MosaicML describes the MPT-7B training dataset as a “MosaicML-curated mix of sources ... [that]  
19 emphasizes English natural language text ... and includes elements of the recently-released RedPajama  
20 dataset.”

21           25.     In a table describing the composition of the MPT-7B training dataset, MosaicML notes  
22 that a large quantity of data comes from a component dataset called “RedPajama—Books”.  
23 MosaicML’s blog post does not further describe the contents of the “RedPajama—Books” dataset.

24           26.     But information about RedPajama is available elsewhere. The RedPajama dataset is  
25 hosted on a website called Hugging Face.<sup>2</sup> According to the documentation for the RedPajama dataset  
26

27 \_\_\_\_\_  
28 <sup>1</sup> Available at <https://www.databricks.com/blog/mpt-7b>

<sup>2</sup> Available at <https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T>

1 found there, its “Books” component is actually a copy of the “Books3 dataset” that is “downloaded  
2 from Huggingface [sic]” when a user runs the script that automatically assembles the RedPajama  
3 dataset. Therefore, anyone who uses the “RedPajama—Books” dataset for training an AI model is  
4 actually using a copy of the Books3 dataset. The documentation for the RedPajama dataset does not  
5 further describe the contents of Books3.

6 27. But information about Books3 is available elsewhere. Books3 is a component of a  
7 separate AI training dataset called The Pile, curated by a research organization called EleutherAI. In  
8 December 2020, EleutherAI introduced this dataset in a paper called “The Pile: An 800GB Dataset of  
9 Diverse Text for Language Modeling.”<sup>3</sup> This paper describes the contents of Books3:

10  
11 Books3 is a dataset of books derived from a copy of the contents of the  
12 Bibliotik private tracker ... Bibliotik consists of a mix of fiction and  
13 nonfiction books and is almost an order of magnitude larger than our  
14 next largest book dataset (BookCorpus2). We included Bibliotik because  
15 books are invaluable for long-range context modeling research and  
16 coherent storytelling.<sup>4</sup>

17 28. Bibliotik is one of a number of notorious “shadow library” websites that also includes  
18 Library Genesis (aka LibGen), Z-Library (aka B-ok), Sci-Hub, and Anna’s Archive. These shadow  
19 libraries have long been of interest to the AI-training community because they host and distribute vast  
20 quantities of unlicensed copyrighted material. For that reason, these shadow libraries also violate the  
21 U.S. Copyright Act.

22 29. The person who assembled the Books3 dataset, Shawn Presser, has confirmed in public  
23 statements that it represents “all of Bibliotik” and contains approximately 196,640 books.

24 30. Plaintiffs’ copyrighted books listed in Exhibit A are among the works in the Books3  
25 dataset. Below, these books are referred to as the **Infringed Works**.

26  
27  
28 <sup>3</sup> Available at <https://arxiv.org/pdf/2101.00027.pdf>

<sup>4</sup> *Id.* at 3–4.

1           31.       Until October 2023, the Books3 dataset was available from Hugging Face. At that time,  
2 the Books3 dataset was removed with a message that it “is defunct and no longer accessible due to  
3 reported copyright infringement.”<sup>5</sup>

4           32.       But before then, anyone who used the “RedPajama—Books” dataset for training  
5 necessarily made a copy of the Books3 dataset. This includes MosaicML, which completed the training  
6 of MPT-7B before May 2023.

7           33.       In sum, MosaicML has admitted training its MPT-7B model on a copy of the  
8 “RedPajama—Books” dataset, which in turn is a copy of the Books3 dataset. Therefore, MosaicML  
9 necessarily trained its MPT-7B model on a copy of Books3. Certain books written by Plaintiffs are part  
10 of Books3—including the Infringed Works—and thus MosaicML necessarily trained MPT-7B on one  
11 or more copies of the Infringed Works, thereby directly infringing the copyrights of the Plaintiffs.

12           34.       In June 2023, MosaicML released another member of the MPT series of large language  
13 models, called MPT-30B. As the name suggests, MPT-30B contained 30 billion weights—over  
14 quadruple the size of MPT-7B—derived from its training dataset. In a table describing the composition  
15 of the MPT-30B training dataset, MosaicML admitted that once again, a large quantity of training data  
16 came from “RedPajama—Books.”<sup>6</sup> By copying the “RedPajama—Books” dataset, MosaicML  
17 necessarily copied Books3—including the Infringed Works—thereby again directly infringing the  
18 copyrights of the Plaintiffs.

19           35.       On information and belief, since it was acquired by Databricks in July 2023, MosaicML  
20 has continued to make copies of the Infringed Works for LLM training and other commercial purposes.  
21  
22  
23  
24  
25  
26

27 \_\_\_\_\_  
28 <sup>5</sup> See [https://huggingface.co/datasets/the\\_pile\\_books3](https://huggingface.co/datasets/the_pile_books3)

<sup>6</sup> Available at <https://www.mosaicml.com/blog/mpt-30b>

**COUNT 1**

**DIRECT COPYRIGHT INFRINGEMENT (17 U.S.C. § 501)**

**AGAINST MOSAICML**

36. Plaintiffs incorporate by reference the preceding factual allegations.

37. As the owners of the registered copyrights in the Infringed Works, Plaintiffs hold the exclusive rights to those books under 17 U.S.C. § 106.

38. To train the MPT-7B and MPT-30B language models, MosaicML copied the Books3 dataset, which includes the Infringed Works. MosaicML made multiple copies of the Books3 dataset during the training of the MPT-7B and MPT-30B models.

39. On information and belief, MosaicML made further copies of the Books3 dataset or subsets thereof to train other models in the MPT family. For instance, MosaicML released a model called MPT-7B-StoryWriter-65k+ (“the StoryWriter model”), a variant of MPT-7B that MosaicML admits was further trained on “a filtered fiction subset of the [B]ooks3 dataset”.<sup>7</sup> The stated purpose of the StoryWriter model is “to read and write stories”—or, put another way, to generate works that directly compete with works in the training dataset.

40. Plaintiffs and the Class members never authorized MosaicML to make copies of their Infringed Works, make derivative works, publicly display copies (or derivative works), or distribute copies (or derivative works). All those rights belong exclusively to Plaintiffs and the Class members under the U.S. Copyright Act.

41. MosaicML repeatedly copied the Infringed Works without Plaintiffs’ permission. MosaicML made these copies without Plaintiffs’ permission and in violation of their exclusive rights under the Copyright Act.

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

---

<sup>7</sup> See <https://www.databricks.com/blog/mpt-7b>

**COUNT 2**  
**VICARIOUS COPYRIGHT INFRINGEMENT**  
**AGAINST DATABRICKS**

43. Plaintiffs incorporate by reference the preceding factual allegations.

44. Databricks acquired MosaicML in July 2023. As the corporate parent of MosaicML, Databricks had the right and ability to control the direct infringements alleged in Count 1 committed by MosaicML, at minimum those occurring after the acquisition. Databricks failed to exert its right and ability to control MosaicML’s infringements.

45. Databricks has directly benefitted financially from the direct infringement alleged in Count 1 because MosaicML generates revenue from its infringing activities, and this revenue belongs to Databricks.

46. Plaintiffs have been injured by Databricks’s acts of vicarious copyright infringement. Plaintiffs are entitled to statutory damages, actual damages, restitution of profits, and other remedies provided by law.

**CLASS ALLEGATIONS**

47. The “**Class Period**” as defined in this Complaint begins on at least March 8, 2021 and runs through the present. Because Plaintiffs do not yet know when the unlawful conduct alleged herein began, but believe, on information and belief, that the conduct likely began earlier than March 8, 2021, Plaintiffs reserve the right to amend the Class Period to comport with the facts and evidence uncovered during further investigation or through discovery.

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

**All persons or entities domiciled in the United States that own a United States copyright in any work that was used as training data for the MPT series of language models during the Class Period.**

1           49.    This Class definition excludes:

- 2           a.    Defendants named herein;
- 3           b.    any of the Defendants' co-conspirators;
- 4           c.    any of Defendants' parent companies, subsidiaries, and affiliates;
- 5           d.    any of Defendants' officers, directors, management, employees, subsidiaries,
- 6                 affiliates, or agents;
- 7           e.    all governmental entities; and
- 8           f.    the judges and chambers staff in this case, as well as any members of their
- 9                 immediate families.

10           50.    **Numerosity.** Plaintiffs do not know the exact number of members in the Class. This  
11 information is in the exclusive control of Defendants. On information and belief, there are at least  
12 thousands of members in the Class geographically dispersed throughout the United States. Therefore,  
13 joinder of all members of the Class in the prosecution of this action is impracticable.

14           51.    **Typicality.** Plaintiffs' claims are typical of the claims of other members of the Class  
15 because Plaintiffs and all members of the Class were damaged by the same wrongful conduct of  
16 Defendants as alleged herein, and the relief sought herein is common to all members of the Class.

17           52.    **Adequacy.** Plaintiffs will fairly and adequately represent the interests of the members of  
18 the Class because the Plaintiffs have experienced the same harms as the members of the Class and have  
19 no conflicts with any other members of the Class. Furthermore, Plaintiffs have retained sophisticated  
20 and competent counsel who are experienced in prosecuting federal and state class actions, as well as  
21 other complex litigation.

22           53.    **Commonality and predominance.** Numerous questions of law or fact common to each  
23 Class member arise from Defendants' conduct and predominate over any questions affecting the  
24 members of the Class individually:

- 25           a.    Whether Defendants violated the copyrights of Plaintiffs and the Class when they
- 26                 obtained copies of Plaintiffs' Infringed Works and copied the Infringed Works into the
- 27                 dataset used to train the MPT models.

- b. Whether Defendants intended to cause further infringement of the Infringed Works with the MPT models because they have distributed these models as so-called “open source” models and advertised those models as a base from which to build further models.
- c. Whether any affirmative defense excuses Defendants’ conduct.
- d. Whether any statutes of limitation limits the potential for recovery for Plaintiffs and the Class.

54. **Other class considerations.** Defendants have acted on grounds generally applicable to the Class. This class action is superior to alternatives, if any, for the fair and efficient adjudication of this controversy. Prosecuting the claims pleaded herein as a class action will eliminate the possibility of repetitive litigation. There will be no material difficulty in the management of this action as a class action. The prosecution of separate actions by individual Class members would create the risk of inconsistent or varying adjudications, establishing incompatible standards of conduct for Defendant.

### DEMAND FOR JUDGMENT

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

- a) This action may proceed as a class action, with Plaintiffs serving as Class Representatives, and with Plaintiffs’ counsel as Class Counsel.
- b) Judgment in favor of Plaintiffs and the Class and against Defendants.
- c) An award of statutory and other damages under 17 U.S.C. § 504 for violations of the copyrights of Plaintiffs and the Class by Defendants.
- d) Reasonable attorneys’ fees as available under 17 U.S.C. § 505 or other applicable statute.
- e) Destruction or other reasonable disposition of all copies Defendants made or used in violation of the exclusive rights of Plaintiffs and the Class, under 17 U.S.C. § 503(b).
- f) Pre- and post-judgment interest on the damages awarded to Plaintiffs and the Class, and that such interest be awarded at the highest legal rate from and after the date this class action complaint is first served on Defendants.

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

4 h) Further relief for Plaintiffs and the Class as may be just and proper.  
5

6 **JURY TRIAL DEMANDED**

7 Under Federal Rule of Civil Procedure 38(b), Plaintiffs demand a trial by jury of all the claims  
8 asserted in this Complaint so triable.  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

1 Dated: March 8, 2024

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

3 Joseph R. Saveri (State Bar No. 130064)  
4 Cadio Zirpoli (State Bar No. 179108)  
5 Christopher K. L. Young (State Bar No. 318371)  
6 Elissa Buchanan (State Bar No. 249996)  
7 **JOSEPH SAVERI LAW FIRM, LLP**  
8 601 California Street, Suite 1000  
9 San Francisco, CA 94108  
10 Telephone: (415) 500-6800  
11 Facsimile: (415) 395-9940  
12 Email: jsaveri@saverilawfirm.com  
13 czirpoli@saverilawfirm.com  
14 cyoung@saverilawfirm.com  
15 eabuchanan@saverilawfirm.com

16 Matthew Butterick (State Bar No. 250953)  
17 1920 Hillhurst Avenue, #406  
18 Los Angeles, CA 90027  
19 Telephone: (323) 968-2632  
20 Facsimile: (415) 395-9940  
21 Email: mb@buttericklaw.com

22 Brian D. Clark (*pro hac vice* pending)  
23 Laura M. Matson (*pro hac vice* pending)  
24 Arielle S. Wagner (*pro hac vice* pending)  
25 Eura Chang (*pro hac vice* pending)  
26 **LOCKRIDGE GRINDAL NAUEN PLLP**  
27 100 Washington Avenue South, Suite 2200  
28 Minneapolis, MN 55401  
Telephone: (612) 339-6900  
Facsimile: (612) 339-0981  
Email: bdclark@locklaw.com  
lmmatson@locklaw.com  
aswagner@locklaw.com  
echang@locklaw.com

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