

[Counsel on signature page]

**UNITED STATES DISTRICT COURT  
NORTHERN DISTRICT OF CALIFORNIA  
SAN FRANCISCO DIVISION**

*In Re Mosaic LLM Litigation*

Master File Case No. 3:24-cv-01451-CRB  
Consolidated with Case No. 3:24-cv-02653-CRB

**SECOND AMENDED CONSOLIDATED  
COMPLAINT**

## Class Action

## Jury Trial Demanded

**REDACTED**

### TABLE OF CONTENTS

I.	OVERVIEW.....	2
II.	JURISDICTION AND VENUE .....	3
III.	PLAINTIFFS .....	3
IV.	DEFENDANTS .....	4
V.	AGENTS AND CO-CONSPIRATORS.....	4
VI.	FACTUAL ALLEGATIONS .....	4
A.	████████ and “RedPajama” Contain Plaintiffs’ and Class Members’ Works .....	4
B.	AI Models Require Training Datasets.....	6
C.	MosaicML Creates a Library of Pirated Books .....	7
D.	Databricks ██████████ .....	9
VII.	CLASS ALLEGATIONS .....	11
VIII.	DEMAND FOR JUDGMENT.....	18
IX.	JURY TRIAL DEMANDED .....	18

19 Plaintiffs Stewart O’Nan, Abdi Nazemian, Brian Keene, Rebecca Makkai, and Jason Reynolds  
20 (together, “Plaintiffs”), on behalf of themselves and all others similarly situated, bring this class-action  
21 complaint (“Complaint”) against defendants Mosaic ML, Inc. (“MosaicML”) and Databricks, Inc.  
22 (“Databricks”) (together “Defendants”).

### I. OVERVIEW

23 1. Defendant MosaicML downloaded hundreds of thousands of copyrighted books without  
24 permission in its quest to develop large language models (“LLMs”). One such model, known as  
25 “Storywriter,” ██████████. This expansive corpus of books—  
26 including ██████████ and “RedPajama”—would form part of a library that  
27

1 MosaicML made [REDACTED]  
2 [REDACTED].

3 2. When Defendant Databricks acquired MosaicML in 2023, MosaicML [REDACTED]  
4 [REDACTED]  
5 [REDACTED]  
6 [REDACTED]  
7 [REDACTED]  
8 [REDACTED]  
9 [REDACTED]

10 3. Plaintiffs and Class members are authors. They own registered copyrights in certain books  
11 that were included in the RedPajama and [REDACTED] datasets that [REDACTED]  
12 [REDACTED] without their permission or compensation.

13 4. Plaintiffs and Class members never authorized Defendants to [REDACTED]  
14 [REDACTED] their copyrighted works. Defendants have never compensated Plaintiffs and  
15 Class members for [REDACTED] their  
16 copyrighted works.

17 5. Defendants also [REDACTED]  
18 [REDACTED]  
19 [REDACTED], and thereby contributed to and induced the infringing activity of others.

20 6. MosaicML benefitted commercially from its acts of copyright infringement, including by  
21 securing investments and contracts with customers for use of its LLMs, and attracting its purchase by  
22 Databricks. Databricks, as the corporate parent of MosaicML and distributor of the MPT and DBRX  
23 models, has also commercially benefitted from this massive copyright infringement by securing  
24 investments and contracts with customers for use of its LLMs.

25 7. Through the above acts, Defendants have jointly infringed Plaintiffs' copyrighted works,  
26 and continue to do so by [REDACTED] the datasets  
27 containing copies of Plaintiffs' and the putative Class's copyrighted books.

## II. JURISDICTION AND VENUE

8. 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).

9. Jurisdiction and venue are proper in this judicial district under 28 U.S.C. § 1391(c)(2) because Defendants are headquartered in this district. MosaicML created the MPT language models and Databricks created the DBRX language models. Defendants distribute these models 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.

10. 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.

### III. PLAINTIFFS

11. Plaintiff Stewart O’Nan is an author who lives in Massachusetts. Mr. O’Nan owns registered copyrights in multiple books, including *Last Night at the Lobster*; *The Night Country*; *Emily, Alone*; *The Circus Fire*; *The Good Wife*; *The Speed Queen*; *The Vietnam Reader*; *West of Sunset*; *City of Secrets*; *Henry, Himself*; and *In the Walled City*.

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

13. Plaintiff Brian Keene is an author who lives in Pennsylvania. Mr. Keene owns registered copyrights in multiple books, including *Ghost Walk*; *City of the Dead*; and *Dead Sea*.

14. Plaintiff Rebecca Makkai is an author who lives in Illinois. Ms. Makkai owns registered copyrights in multiple books, including, *The Hundred Year House*; *Music for Wartime*; and *The Great Believers*.

15. Plaintiff Jason Reynolds is an author who lives in Washington, D.C. Mr. Reynolds owns registered copyrights in multiple books, including *As Brave as You*; *When I was the Greatest*; *All American Boys*; *Ghost*; *Patina*; *Long Way Down*; *Sunny*; *For Every One*; and *Look Both Ways*.

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

## IV. DEFENDANTS

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

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

## V. AGENTS AND CO-CONSPIRATORS

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

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

## VI. FACTUAL ALLEGATIONS

## A. **“RedPajama” and “RedPajama.com” and “RedPajama” and “RedPajama” Contain Plaintiffs’ and Class Members’ Works**

21. The Pile is a dataset curated by a research organization called EleutherAI for use in training AI models. In December 2020, EleutherAI introduced this dataset in a paper called “The Pile:

1 An 800GB Dataset of Diverse Text for Language Modeling.”<sup>1</sup> The paper provides a description of  
 2 Books3, a dataset contained within The Pile:

3 Books3 is a dataset of books derived from a copy of the contents of the Bibliotik  
 4 private tracker ... Bibliotik consists of a mix of fiction and nonfiction books and  
 5 is almost an order of magnitude larger than our next largest book dataset  
 (BookCorpus2). We included Bibliotik because books are invaluable for long-  
 range context modeling research and coherent storytelling.<sup>2</sup>

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

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

13 24. Until October 2023, the Books3 dataset was available from Hugging Face as a standalone  
 14 dataset. At that time, the Books3 dataset was removed with a message that it “is defunct and no longer  
 15 accessible due to reported copyright infringement.”<sup>3</sup>

16 25. Books3 was included within another dataset known as RedPajama. Released in April 2023  
 17 and created by the company Together AI, the RedPajama dataset contained a subset called “Books” (also  
 18 referred to as RedPajama-Books) that was actually a copy of the “Books3 dataset” that was “downloaded  
 19 from Huggingface [sic].” A user could either download the dataset or run scripts that automatically  
 20 assembled the RedPajama dataset.<sup>4</sup> After the “Books3 dataset” was removed from Hugging Face in  
 21  
 22

---

23 <sup>1</sup> Available at Gao, Leo, et al., *The Pile: An 800GB Dataset of Diverse Text for Language Modeling*,  
 24 arXiv:2101.00027 (Dec. 31, 2020), <https://arxiv.org/pdf/2101.00027.pdf>.

25 <sup>2</sup> *Id.* at 3–4.

26 <sup>3</sup> See *The Pile Books3 Dataset*, Hugging Face, [https://web.archive.org/web/20231127101818/https://huggingface.co/datasets/the\\_pile\\_books3](https://web.archive.org/web/20231127101818/https://huggingface.co/datasets/the_pile_books3) (on file  
 27 with the Joseph Saveri Law Firm, LLP) (last visited Nov. 18, 2025).

28 <sup>4</sup> Available at TogetherComputer, *RedPajama-Data-1T Dataset*, Hugging Face, <https://web.archive.org/web/20230420075601/https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T> (on file with the Joseph Saveri Law Firm, LLP) (last visited Nov. 18, 2025).

1 October 2023, the RedPajama dataset documentation also added a message that Books3 is defunct “due  
 2 to reported copyright infringement.”<sup>5</sup>

3 26. But before October 2023, anyone who downloaded the “RedPajama” dataset was  
 4 downloading a copy of the Books3 dataset.

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

7 **B. AI Models Require Training Datasets**

8 28. A large language model (“LLM”) is AI software designed to emit convincingly  
 9 naturalistic text outputs in response to user prompts.

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

14 30. Training consists of a multi-stage process (known as the training pipeline) that includes  
 15 the acquisition and curation of the dataset, processing of the dataset, feeding the dataset into the model  
 16 so that the model can extract the patterns and relationships from the protected expression contained  
 17 therein, and further fine-tuning the model for more specialized uses with even more data.

18 31. The first step in training the model is acquiring and curating the data that goes in to the  
 19 model. Training an LLM is not only a function of quantity of data, but also of quality. The selection and  
 20 curation of training data is therefore an important first step in training. Copyrighted books tend to be  
 21 high-quality data for training LLMs.

22 32. During training, the LLM copies and ingests each textual work in the training dataset and  
 23 extracts protected expression from it. During what is known as *pretraining*, the LLM progressively  
 24 adjusts its output to more closely approximate the protected expression copied from the training dataset.  
 25 The LLM records the results of this process in a large set of numbers called weights (also known as  
 26

---

27 <sup>5</sup> Available at TogetherComputer, *RedPajama-Data-1T Dataset*, Hugging Face,  
 28 <https://web.archive.org/web/20240510231649/https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T> (on file with the Joseph Saveri Law Firm, LLP) (last visited Nov. 18, 2025).

1 *parameters*) that are stored within the model. These weights are entirely and uniquely derived from the  
 2 protected expression in the training dataset. Once a model is pretrained, it results in a trained model  
 3 known as a *base* or *foundational* model.

4       33. During the development process of an LLM, engineers may also conduct experiments  
 5 known as “ablation studies” that test the effect of certain data on the model. This can include, for  
 6 example, determining whether there is a difference in the quality of a model’s output if it is trained with  
 7 books versus without. A dataset may be used to run such experiments but be excluded from the final  
 8 training dataset of the model. Importantly, these datasets too may consist of copyrighted works,  
 9 including books.

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

15      35. Throughout each step of the training pipeline, the same dataset may be used multiple  
 16 times. Indeed, given the cost of developing an LLM, it is a ubiquitous practice to retain datasets for  
 17 future use, whether that use is to pretrain other models, to perform ablations on a model, or to fine-tune  
 18 an already trained base model. The implication is that if a dataset contains unlawfully-obtained  
 19 copyrighted material, each step of the training pipeline may result in an unauthorized use (i.e.,  
 20 infringement) of that copyrighted work.

21      **C. MosaicML Creates a Library of Pirated Books**

22      36. MosaicML was founded in 2020 to provide tools to facilitate the training of AI models.

23      37. One of these tools is the MosaicML Platform, which companies or individuals can use to  
 24 train their own AI models.<sup>6</sup> The MosaicML Platform allows users to select their training data, including  
 25 through a feature known as MosaicML Streaming that allows customers to access datasets stored on  
 26

---

27  
 28      <sup>6</sup> Available at Mosaic ML, Inc. (organization page), <https://huggingface.co/mosaicml>, (on file with the  
 Joseph Saveri Law Firm, LLP) (last visited Nov. 20, 2025).

1 remote or cloud servers.<sup>7</sup> MosaicML provides its customers with selected datasets to use with the  
 2 MosaicML Platform, [REDACTED].

3 38. [REDACTED]

4 [REDACTED]

5 [REDACTED].

6 39. In April 2023, MosaicML obtained a copy of the RedPajama dataset to use to develop its  
 7 MPT models. RedPajama also contains a subset called "RedPajama-Books" that includes the Books3  
 8 dataset containing Plaintiffs' and Class members' works.

9 40. MosaicML then [REDACTED] and

10 RedPajama [REDACTED]

11 [REDACTED]

12 They contain close to [REDACTED] books combined.

13 41. MosaicML used this trove of pirated books for its own benefit. [REDACTED]

14 [REDACTED]

15 [REDACTED]

16 [REDACTED]

17 [REDACTED], MosaicML Chief Technology Officer Hanlin Tang [REDACTED]

18 [REDACTED]

19 42. At around the same time, in April 2023, MosaicML [REDACTED]

20 [REDACTED]. MosaicML's then-Chief Scientist Jonathan Frankle [REDACTED]

21 [REDACTED] Naveen  
 22 Rao, former CEO of MosaicML and Vice President of AI at Databricks, stated [REDACTED]

23 [REDACTED]

24 [REDACTED]

25 [REDACTED]

26 <sup>7</sup> *Id.*

27 8 [REDACTED],

28 [REDACTED] (Apr.

19,2023).

1       43.     In May 2023, MosaicML released the first in its MPT series of large language models,  
 2 called MPT-7B.

3       44.     In a blog post called “Introducing MPT-7B: A New Standard for Open-Source,  
 4 Commercially Usable LLMs,” MosaicML describes the MPT-7B training dataset as a “MosaicML-  
 5 curated mix of sources . . . [that] emphasizes English natural language text . . . and includes elements of  
 6 the recently-released RedPajama dataset.”<sup>9</sup>

7       45.     In a table describing the composition of the MPT-7B training dataset, MosaicML notes  
 8 that a large quantity of that training data came from “RedPajama—Books”.

9       46.     MosaicML then released a model called MPT-7B-StoryWriter-65k+ (“the Storywriter  
 10 model”), a variant of MPT-7B that MosaicML admits was further trained on “a filtered fiction subset of  
 11 the [B]ooks3 dataset.” The stated purpose of the Storywriter model is “to read and write stories”—or,  
 12 put another way, to generate works that directly compete with works in the training dataset. A MosaicML  
 13 employee highlighted that [REDACTED]

14 [REDACTED]  
 15       47.     In June 2023, MosaicML released another member of the MPT series of large language  
 16 models, called MPT-30B. As the name suggests, MPT-30B contained 30 billion weights—over  
 17 quadruple the size of MPT-7B—derived from its training dataset. In a table describing the composition  
 18 of the MPT-30B training dataset, MosaicML admitted that once again, a large quantity of that training  
 19 data came from “RedPajama—Books.”<sup>10</sup>

20       **D. Databricks** [REDACTED]

21       48.     After Databricks acquired MosaicML in July 2023, Databricks took [REDACTED]

22 [REDACTED]  
 23 [REDACTED]  
 24  
 25  
 26  
 27 <sup>9</sup> See MosaicML, *Introducing MPT-7B: A New Standard for Open-Source, Commercially Usable LLMs*, Databricks Blog (May 5, 2023), <https://www.databricks.com/blog/mpt-7b>.

28 <sup>10</sup> Available at Mosaic Research, Mosaic Research Blog, <https://www.mosaicml.com/blog/mpt-30b>. (on file with the Joseph Saveri Law Firm, LLP) (last visited Nov. 20, 2025).

1 49. Databricks [REDACTED]  
2 [REDACTED]  
3 [REDACTED]  
4 50. Databricks [REDACTED]  
5 [REDACTED]. Databricks also continued to  
6 operate the MosaicML Platform that enabled clients to train their own AI models [REDACTED]  
7 [REDACTED]<sup>11</sup> Databricks [REDACTED]  
8 [REDACTED]  
9 [REDACTED]. In doing so, Databricks [REDACTED]  
10 [REDACTED]  
11 51. [REDACTED]  
12 [REDACTED]  
13 [REDACTED].  
14 52. Soon after acquiring MosaicML, Databricks [REDACTED]  
15 [REDACTED]. Referred to as [REDACTED], Databricks employees [REDACTED]  
16 [REDACTED]. A Databricks employee [REDACTED]  
17 [REDACTED]  
18 [REDACTED] Databricks would later [REDACTED].  
19 53. Plainly, [REDACTED] Databricks' acquisition of  
20 MosaicML. As Chief Scientist Jonathan Frankle stated [REDACTED]  
21 [REDACTED] And Naveen Rao stated that  
22 Databricks [REDACTED]  
23 [REDACTED].  
24  
25  
26  
27

---

28<sup>11</sup> See Mosaic ML, Inc. (organization page), <https://huggingface.co/mosaicml> (on file with the Joseph Saveri Law Firm, LLP) (last visited Nov. 20, 2025).

54. On March 27, 2024, Databricks announced its DBRX series of LLMs in a blog post called “Introducing DBRX: A New State-of-the-Art Open LLM.”<sup>12</sup> In the blog post, Databricks described DBRX as a “general-purpose LLM,” that was trained on a “curated” dataset of 12 trillion tokens. There are two versions of DBRX—“DBRX Base” and “DBRX Instruct.”

55. Defendants have not publicly identified the training data used for these DBRX models. Databricks' Vice President of AIs and former CEO of MosaicML, however, has stated that DBRX was trained on "open data sets that the community knows."<sup>13</sup>

56. Defendants have publicly stated that the DBRX models were the culmination of “years of LLM development at Databricks that includes the MPT . . . projects” and that “[t]he development of DBRX was led by the Mosaic team that previously built the MPT model family.”<sup>14</sup>

57. There was nothing fair about [REDACTED]  
[REDACTED] MosaicML downloaded these books [REDACTED]  
and RedPajama—without permission from or compensation to their authors. [REDACTED]

## VII. CLASS ALLEGATIONS

58. 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.

<sup>12</sup> Available at Mosaic Research Team, *Introducing DBRX: A New State-of-the-Art Open LLM*, Mosaic AI Research (Mar. 27, 2024), <https://www.databricks.com/blog/introducing-dbrx-new-state-art-open-llm>.

<sup>13</sup> Available at Kyle Wiggers, *Databricks Spent \$10M on a Generative AI Model That Still Can't Beat GPT-4*, TechCrunch (Mar. 27, 2024), <https://techcrunch.com/2024/03/27/databricks-spent-10m-on-a-generative-ai-model-that-still-cant-beat-gpt-4/>.

<sup>14</sup> Available at Mosaic Research Team, *Introducing DBRX: A New State-of-the-Art Open LLM*, Mosaic AI Research (Mar. 27, 2024), <https://www.databricks.com/blog/introducing-dbrx-new-state-art-open-llm>.

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

4       **All legal or beneficial owners of copyrighted works that: (A) were registered with the  
 5       United States Copyright Office within five years of the work's first publication; (B) were  
 6       downloaded, reproduced, or distributed by Defendants; (C) were registered with the United  
 7       States Copyright Office before being downloaded, reproduced, or distributed by  
 8       Defendants, or were registered within three months of first publication; and (D) are  
 9       assigned one or more International Standard Books Number(s) (ISBN) or Amazon Standard  
 10       Identification Number(s) (ASIN).**

11       60.   These Class definitions exclude:

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

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

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

25       63.   **Adequacy.** Plaintiffs will fairly and adequately represent the interests of the members of  
 26 the Class because the Plaintiffs have experienced the same harms as the members of the Class and have  
 27 no conflicts with any other members of the Class. Furthermore, Plaintiffs have retained sophisticated and

competent counsel who are experienced in prosecuting federal and state class actions, as well as other complex litigation.

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

a. Whether Defendants violated the copyrights of Plaintiffs and the Class when they

b. Whether Defendants intended to cause further infringement of the Infringed Works by virtue of [REDACTED]  
[REDACTED]  
[REDACTED].

c. Whether any affirmative defense excuses Defendants' conduct.

d. Whether any statutes of limitation limits the potential for recovery for Plaintiffs and the Class.

65. **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.

## COUNT 1

## **Direct Copyright Infringement (17 U.S.C. § 501) against all Defendants**

66. Plaintiffs incorporate by reference the preceding factual allegations.

67. 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.

1 68. MosaicML downloaded copies of the RedPajama and [REDACTED] datasets [REDACTED]  
2 [REDACTED], which includes the Infringed Works.

3 69. MosaicML stored copies of RedPajama [REDACTED]  
4 [REDACTED]  
5 [REDACTED].

6 70. MosaicML [REDACTED]  
7 [REDACTED].

8 71. MosaicML also [REDACTED]  
9 [REDACTED]  
10 [REDACTED]. Upon information and belief, [REDACTED]

11 [REDACTED]  
12 72. To develop the MPT-7B, MPT-30B, and MPT-7B-StoryWriter-65k+ language models,  
13 MosaicML copied [REDACTED] and RedPajama [REDACTED], and train  
14 these models. MosaicML engineers further [REDACTED]  
15 [REDACTED]. MosaicML made multiple copies [REDACTED]  
16 [REDACTED].

17 73. When Databricks acquired MosaicML in July 2023, [REDACTED]  
18 [REDACTED]  
19 [REDACTED]. MosaicML executives became Databricks executives and ran the pre-existing  
20 MosaicML business as part of Databricks. After the acquisition, Databricks [REDACTED]  
21 [REDACTED]  
22 [REDACTED].

23 74. Databricks further [REDACTED]  
24 [REDACTED].

25 75. Databricks marketed to its customers the ability to use the MosaicML Platform to train  
26 LLM models [REDACTED]  
27 [REDACTED]  
28 [REDACTED].

1 76. [REDACTED], Databricks [REDACTED]  
2 [REDACTED]  
3 [REDACTED].

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

8 78. By [REDACTED] datasets containing  
9 copies of Plaintiffs' Works, Defendants directly infringed Plaintiffs' exclusive rights in their copyrighted  
10 works.

11 79. By copying, storing, processing, and reproducing the MPT models trained on Plaintiffs'  
12 Works, MosaicML has directly infringed Plaintiffs' exclusive rights in their copyrighted works.

13 80. Defendants' infringing conduct alleged herein was and continues to be willful and carried  
14 out with full knowledge of Plaintiffs' rights in the copyrighted works. As a direct result of their conduct,  
15 Defendants have wrongfully profited from copyrighted works that they do not own.

16 81. By and through the actions alleged above, Defendants have infringed and will continue to  
17 infringe Plaintiffs' copyrights.

18 82. Plaintiffs have been and continue to be injured by Defendants' acts of direct copyright  
19 infringement. Plaintiffs are entitled to statutory damages, actual damages, restitution of profits, and all  
20 appropriate legal and equitable relief.

21 **COUNT 2**

22 **Vicarious Copyright Infringement against Databricks**

23 83. Plaintiffs incorporate by reference the preceding factual allegations.

24 84. Databricks acquired MosaicML in July 2023. As the corporate parent of MosaicML,  
25 Databricks had the right and ability to control the direct infringements alleged in Count 1 committed by  
26 MosaicML, at minimum those occurring after the acquisition. Databricks failed to exercise its right and  
27 ability to control MosaicML's infringements.

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

86. Plaintiffs have been and continue to be injured by Databricks's acts of vicarious copyright infringement. Plaintiffs are entitled to statutory damages, actual damages, restitution of profits, and all appropriate legal and equitable.

## COUNT 3

## Contributory Copyright Infringement against all Defendants

87. Plaintiffs incorporate by reference the preceding factual allegations.

88. Defendants [REDACTED]

## Defendants

89. 

violates Plaintiffs and Class members' exclusive rights under 17 U.S.C. § 106.

90. Defendants [REDACTED] Such use of the Infringed Works

91. Defendants are contributorily liable for these direct infringements by [REDACTED],

92. Defendants are well aware of [REDACTED] infringing activity. Defendants facilitated, encouraged, and materially contributed to such infringement, including but not limited to by

10 of 10 | Page | Page Number | Page Total | Page Footer

93. Defendants failed to take steps to stop the specific infringing activity. Defendants [REDACTED]

[REDACTED] . As a direct and

1 proximate result, [REDACTED] have infringed  
 2 Plaintiffs and Class members' copyrights in their works.

3 94. By [REDACTED] Plaintiff and Class members Works, Defendants  
 4 violated the exclusive rights under 17 U.S.C. § 106 of Plaintiffs and Class members.

5 95. Plaintiffs have been and continue to be injured by Defendants' acts of contributory  
 6 copyright infringement. Plaintiffs are entitled to statutory damages, actual damages, restitution of profits,  
 7 and all appropriate legal and equitable relief.

8 **COUNT 4**

9 **Inducement of Copyright Infringement against all Defendants**

10 96. Plaintiffs incorporate by reference the preceding factual allegations.

11 97. Defendants [REDACTED]

12 [REDACTED]  
 13 [REDACTED] Defendants did so by [REDACTED]

14 [REDACTED].

15 [REDACTED].

16 98. At least [REDACTED]

17 [REDACTED]  
 18 [REDACTED]  
 19 [REDACTED] violates Plaintiffs and Class members' exclusive rights under 17 U.S.C. § 106.

20 99. Defendants made a material contribution to this infringing activity by [REDACTED]

21 [REDACTED]  
 22 [REDACTED].

23 100. Plaintiffs have been and continued to be injured by Defendants' acts of inducement of  
 24 copyright infringement. Plaintiffs are entitled to statutory damages, actual damages, restitution of profits,  
 25 and all appropriate legal and equitable relief.

26 [REDACTED]

27 [REDACTED]

28 [REDACTED]

## **VIII. 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 and reimbursement of costs under 17 U.S.C. § 505 or otherwise.
- e) A declaration that such infringement is willful.
- f) 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).
- g) 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.
- h) Defendants are to be jointly and severally responsible financially for the costs and expenses of a Court-approved notice program through post and media designed to give immediate notification to the Class.
- i) Further relief for Plaintiffs and the Class as may be appropriate.

## **IX. JURY TRIAL DEMANDED**

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

1 Dated: January 21, 2026

Respectfully submitted,

2 By: /s/ Joseph R. Saveri

3 Joseph R. Saveri (SBN 130064)  
 4 Christopher K.L. Young (SBN 318371)  
 5 Evan Creutz (SBN 349728)  
 6 Elissa A. Buchanan (SBN 249996)  
 7 William Castillo Guardado (SBN 294159)  
 Holden Benon (SBN 325847)  
**JOSEPH SAVERI LAW FIRM, LLP**  
 601 California Street, Suite 1505  
 San Francisco, CA 94108  
 Telephone: (415) 500-6800  
 Facsimile: (415) 395-9940  
 jsaveri@saverilawfirm.com  
 cyoung@saverilawfirm.com  
 ecreutz@saverilawfirm.com  
 eabuchanan@saverilawfirm.com  
 wcastillo@saverilawfirm.com  
 hbenon@saverilawfirm.com

13  
 14 Matthew Butterick (SBN 250953)  
 15 1920 Hillhurst Avenue, #406  
 16 Los Angeles, CA 90027  
 17 Telephone: (323) 968-2632  
 18 Facsimile: (415) 395-9940  
 19 mb@buttericklaw.com

20 Bryan L. Clobes (admitted *pro hac vice*)  
 21 Mohammed A. Rathur (admitted *pro hac vice*)  
 22 Nabihah Maqbool (admitted *pro hac vice*)  
**CAFFERTY CLOBES MERIWETHER**  
**& SPRENGEL LLP**  
 23 135 South LaSalle Street, Suite 3210  
 Chicago, IL 60603  
 Telephone: (312) 782-4880  
 bclobes@caffertyclobes.com  
 mrathur@caffertyclobes.com  
 nmaqbool@caffertyclobes.com

Justin A. Nelson (admitted *pro hac vice*)  
Alejandra C. Salinas (admitted *pro hac vice*)  
**SUSMAN GODFREY L.L.P.**  
1000 Louisiana Street, Suite 5100  
Houston, TX 77002-5096  
Telephone: (713) 651-9366  
[jnelson@susmangodfrey.com](mailto:jnelson@susmangodfrey.com)  
[asalinas@susmangodfrey.com](mailto:asalinas@susmangodfrey.com)

Rohit D. Nath (SBN 316062)  
**SUSMAN GODFREY L.L.P**  
1900 Avenue of the Stars, Suite 1400  
Los Angeles, CA 90067-2906  
Telephone: (310) 789-3100  
RNath@susmangodfrey.com

Elisha Barron (admitted *pro hac vice*)  
Craig Smyser (admitted *pro hac vice*)  
**SUSMAN GODFREY L.L.P**  
One Manhattan West, 51st Floor  
New York, NY 10019  
Telephone: (212) 336-8330  
[ebarron@susmangodfrey.com](mailto:ebarron@susmangodfrey.com)  
[csmysyer@susmangodfrey.com](mailto:csmysyer@susmangodfrey.com)

Jordan W. Connors (admitted *pro hac vice*)  
Trevor D. Nystrom (admitted *pro hac vice*)  
Dylan Salzman (admitted *pro hac vice*)  
**SUSMAN GODFREY L.L.P**  
401 Union Street, Suite 3000  
Seattle, WA 98101  
Telephone: (206) 516-3880  
[jconnors@susmangodfrey.com](mailto:jconnors@susmangodfrey.com)  
[tnystrom@susmangodfrey.com](mailto:tnystrom@susmangodfrey.com)  
[dsalzman@susmangodfrey.com](mailto:dsalzman@susmangodfrey.com)

Rachel J. Geman (admitted *pro hac vice*)  
Danna Z. Elmasry (admitted *pro hac vice*)  
**LIEFF CABRASER HEIMANN &**  
**BERNSTEIN, LLP**  
250 Hudson Street, 8th Floor  
New York, NY 10013  
Telephone: (212) 355-9500  
[rgeman@lchb.com](mailto:rgeman@lchb.com)  
[delmasry@lchb.com](mailto:delmasry@lchb.com)

1 Anne B. Shaver  
2 **LIEFF CABRASER HEIMANN &**  
3 **BERNSTEIN, LLP**  
4 275 Battery Street, 29th Floor  
5 San Francisco, CA 94111  
6 Telephone: (415) 956-1000  
7 ashaver@lchb.com

8 Betsy A. Sugar (admitted *pro hac vice*)  
9 Kenneth S. Byrd (admitted *pro hac vice*)  
10 **LIEFF CABRASER HEIMANN &**  
11 **BERNSTEIN, LLP**  
12 222 2nd Avenue S. Suite 1640  
13 Nashville, TN 37201  
14 Telephone: (615) 313-9000  
15 bsugar@lchb.com  
16 kbyrd@lchb.com

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

19  
20  
21  
22  
23  
24  
25  
26  
27  
28