

Eugene Y. Turin (SB # 342413)
MC GUIRE LAW, P.C.
1089 Willowcreek Road, Suite 200
San Diego, CA 92131
Tel: (312) 893-7002 Ex. 3
Fax: 312-275-7895
eturin@mcgpc.com

Counsel for Plaintiff and the Putative Class Members

UNITED STATES DISTRICT COURT
NORTHERN DISTRICT OF CALIFORNIA
SAN FRANCISCO DIVISION

DARIUS H. JAMES, individually and
on behalf of similarly situated
individuals,) Case No.
Plaintiff,) **COMPLAINT**
v.) **CLASS ACTION**
TOGETHER COMPUTER, INC., a) **DEMAND FOR JURY TRIAL**
Delaware corporation,)
Defendant.)

1 Plaintiff Darius H. James (“Plaintiff”), on behalf of himself and all others
 2 similarly situated, bring this class action complaint (“Complaint”) against Defendant
 3 Together Computer, Inc. (“Together AI” or “Defendant”).

4 **OVERVIEW**

5 1. Artificial intelligence (“AI”) refers to software engineered to mimic
 6 human-like reasoning and inference through algorithmic processes, typically
 7 leveraging statistical methods.

8 2. Large language models (“LLM”) are AI software programs designed to
 9 reply to user prompts with natural-sounding text outputs.

10 3. While the traditional coding process involves human coders inputting
 11 explicit instructions, an LLM is instead trained by processing vast quantities of text
 12 from diverse sources (a “training dataset”), learning statistical patterns and
 13 associations within that data, and encoding those abstract representations into a vast
 14 array of numerical values known as parameters. These parameters emerge from the
 15 collective structure of the entire training dataset, which often blends public domain,
 16 licensed, and, crucially, unlicensed copyrighted materials. When an LLM generates
 17 text in response to a prompt, it performs probabilistic computations over these
 18 parameters to produce coherent, contextually relevant output that generalizes from the
 19 learned patterns.

20 4. Plaintiff and Class members are authors. They own registered copyrights
 21 in certain books (the “Infringed Works”) that were included in the open training
 22 dataset, RedPajama, that Together AI assembled and published for the purpose of
 23 training LLMs, and ultimately driving such LLM development to its own
 24 development platform. Plaintiff and Class members never authorized Together AI to
 25 download, copy, store, publish or use their copyrighted works to create the RedPajama
 26 dataset.

27 5. By the acts described in further detail below, Together AI infringed on
 28

1 Plaintiff's copyrighted works, and it may continue to do so if it continues to download,
 2 copy, store, publish and use the RedPajama dataset which contains copies of
 3 Plaintiff's and the putative Class's Infringed Works.

4 **JURISDICTION AND VENUE**

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

7 7. Jurisdiction and venue are proper in this judicial district under 28 U.S.C.
 8 § 1391(c)(2) because Together AI is headquartered in this district. Together AI copied
 9 Plaintiff's and Class members' Infringed Works and assembled the RedPajama
 10 dataset. Therefore, a substantial part of the events giving rise to the claim occurred in
 11 this District. A substantial portion of the affected interstate trade and commerce was
 12 carried out in this District. Defendant has transacted business, maintained substantial
 13 contacts, and/or committed overt acts in furtherance of the illegal scheme and
 14 conspiracy throughout the United States, including in this District. Defendant's
 15 conduct has had the intended and foreseeable effect of causing injury to persons
 16 residing in, located in, or doing business throughout the United States, including in
 17 this District.

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

22 **PARTIES**

23 9. Plaintiff Darius H. James is an author and performance artists who
 24 resides in Connecticut. He is the author of two American publications, *Negrophobia*
 25 and *That's Blaxploitation!: Roots of the Baadasssss 'Tude*, and two bilingual German
 26 publications, *Voodoo Stew* and *Froggie Chocolates' Christmas Eve*.

27 10. A list of Plaintiff's registered copyrights are attached hereto as Exhibit

A.

11. Defendant Together AI is a Delaware corporation with its principal place of business at 251 Rhode Island Street, Suite 205, San Francisco, CA 94103.

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

13. 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, or for Defendant with respect to the acts, violations, and common course of conduct alleged herein.

FACTUAL ALLEGATIONS

14. Together AI is a company that contributes open-source research, models, and datasets to other AI-focused companies and sells access to a cloud platform to help other entities train, fine-tune, and deploy generative AI models.¹

15. In an effort to drive more AI startups to its cloud platform, TogetherAI orchestrated and assembled the RedPajama dataset, a dataset comprised of a mixture of publicly available works as well as copyrighted works.

16. The RedPajama dataset assembled and published by TogetherAI contained within it a deduplicated copy of the Books3 dataset. Books3 was described in a paper by EluetherAI called “*The Pile: An 800GB Dataset of Diverse Text for Language Modeling*” as follows:

¹ <https://www.together.ai/about>

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

6 17. Before October 2023, Books3 was available for download from Hugging
 7 Face (a website dedicated to “a mission to democratize good machine learning”) as a
 8 standalone dataset.³ But in October 2023, the Books3 dataset was removed with a
 9 message that it “is defunct and no longer accessible due to reported copyright
 10 infringement.”⁴

11 18. Plaintiff’s copyrighted books are among the works in the Books3 dataset.

12 19. Books3 was also downloaded and copied by Defendant and utilized to
 13 create and publish its RedPajama dataset in April of 2023.

14 20. The RedPajama dataset contained a subset called “Books” or
 15 “RedPajama-Books” that was actually a copy of the Books3 dataset. Specifically, the
 16 RedPajama dataset “is a publicly available, fully open, best-effort reproduction of the
 17 training data. . . used to train the first iteration of LLaMA family of models.”⁵ This
 18 LLaMA training dataset included the Books3 section of The Pile (a broader publicly
 19 available dataset which contained Books3).⁶

20 21. Thus, RedPajama training datasets, including the RedPajama-Books
 21 subset, contained copyrighted works, including books written and copyrighted by
 22

23 _____
 24 ² <https://arxiv.org/abs/2101.00027>

25 ³ <https://huggingface.co/huggingface>

26 ⁴ https://web.archive.org/web/20231127101818/https://huggingface.co/datasets/the_pile_books3

27 ⁵ <https://arxiv.org/html/2411.12372v1>

28 ⁶ <https://arxiv.org/pdf/2302.13971>

1 Plaintiff and Class members.

2 22. Together AI was chiefly responsible for the creation of the RedPajama
 3 dataset, acknowledging itself as collaborator in its blog post and announcing the
 4 project publicly on the X platform on April 17, 2023.⁷⁸ Importantly, the RedPajama
 5 dataset is hosted on Defendant's specific GitHub organization page
 6 ("togethercomputer")⁹ and its affiliation with the RedPajama dataset page is
 7 evidenced by the license notes disclosing "Copyright 2023 Together Computer."¹⁰

8 23. Defendant's blog post announcing the creation of the RedPajama dataset
 9 also provides a link to download the dataset on its Hugging Face "togethercomputer"
 10 profile page.¹¹

11 24. While the Hugging Face website currently states that Books3 was
 12 removed from RedPajama due to "reported" copyright infringement, an archived
 13 version of that Hugging Face webpage from April 2023 explicitly disclosed that
 14 Books3 was included in the RedPajama corpus of materials.^{12 13}

15 25. On this same archived webpage, it also states that "we [Together AI] use
 16 simhash to remove near duplicates" from Books3, meaning that Together AI
 17 downloaded, copied, deduplicated, repackaged and re-published Books3 within its
 18 RedPajama dataset online for others to download and train their LLMs.¹⁴

19 26. Accordingly, through its creation and publication of the RedPajama
 20 dataset, and other data subsets including RedPajama-Books, Defendant has itself

22 ⁷ <https://www.together.ai/blog/redpajama>

23 ⁸ <https://x.com/togethercompute/status/1647917989264519174>

24 ⁹ <https://github.com/togethercomputer/RedPajama-Data>

¹⁰ *Id.*

¹¹ <https://www.together.ai/blog/redpajama>

¹² <https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T>

¹³ <https://web.archive.org/web/20230417120911/https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T>

¹⁴ *Id.*

1 downloaded, copied, stored, and used Books3 and the copyrighted works of Plaintiff
2 and the members of the Class.

3 27. Furthermore, Defendant has actively promoted the RedPajama dataset
4 and incentivized third parties to utilize, download, and copy it. For example, in May
5 2023, Together AI announced a monetary incentive, hosted by Chai Research,
6 encouraging third parties to “create the best chatbot by fine-tuning RedPajama-
7 INCITE-3B, with a \$1M prize for the winner!”¹⁵

8 28. Between April 2023 and October 2023, Together AI's promotions
9 incentivized over 190,000 downloads of the RedPajama dataset within the Together
10 AI community and third parties.¹⁶

11 29. Thus, from time that Together AI made RedPajama available on the
12 Hugging Face and GitHub websites in April of 2023, to the time that Books3 was
13 removed from the dataset, Defendant made available for download, and promoted the
14 downloading and copying of a vast repository of texts including the copyrighted
15 works of Plaintiff and the members of the Class.

CLASS ALLEGATIONS

17 30. The “Class Period” as defined in this Complaint begins on at least
18 November 5, 2022 and runs through the present. Because Plaintiff does not yet know
19 when the unlawful conduct alleged herein began, but believes, on information and
20 belief, that the conduct likely began prior to November 5, 2022, Plaintiff reserves the
21 right to amend the Class Period to comport with the facts and evidence uncovered
22 during further investigation or through discovery.

¹⁵ <https://www.together.ai/blog/redpajama-3b-updates#:~:text=%20for%20longer%20context>

¹⁶ <https://www.together.ai/blog/redpajama-data-v2#:~:text=Over%20the%20last%20half%20a,released%20specifically%20for%20LLM%20training>

1 31. Plaintiff seeks certification of the following Class pursuant to Federal
2 Rules of Civil Procedure 23(a), 23(b)(2), and 23(b)(3):

3 4 All persons or entities domiciled in the United States that own a United
4 States copyright in any work that was downloaded, copied, stored, used,
5 or contained within the RedPajama training datasets during the Class
6 Period.

7 32. Plaintiff will fairly and adequately represent and protect the interests of
8 the other members of the Class. Plaintiff has retained counsel with substantial
9 experience in prosecuting complex litigation and class actions. Plaintiff and their
10 counsel are committed to vigorously prosecuting this action on behalf of the other
11 members of the Class, and have the financial resources to do so. Neither Plaintiff nor
12 their counsel have any interest adverse to those of the other members of the Class.

13 33. Absent a class action, most members of the Class would find the cost of
14 litigation their claims to be prohibitive and would have no effective remedy. The class
15 treatment of common questions of law and fact is also superior to multiple individual
16 actions or piecemeal litigation in that it conserves the resources of the courts and the
17 litigants and promotes consistency and efficiency of adjudication.

18 34. Defendant has acted and failed to act on grounds generally applicable to
19 Plaintiff and the other members of the Class, requiring the Court's imposition of
20 uniform relief to ensure compatible standards of conduct toward the members of the
21 Class, and making injunctive or corresponding declaratory relief appropriate for the
22 Class as a whole.

23 35. The factual and legal bases of Defendant's liability to Plaintiff and to the
24 other members of the Class are the same, resulting in injury to Plaintiff and to all of
25 the other members of the Class. Plaintiff and the other members of the Class have all
26 suffered harm and damages as a result of Defendant's unlawful and wrongful conduct.

27 36. There are many questions of law and fact common to the claims of

1 Plaintiff and the other members of the Class, and those questions predominate over
2 any questions that may affect individual members of the Class. Common questions
3 for the Class include, but are not limited to, the following:

4 a. Whether Defendant violated the copyrights of Plaintiff and the Class by
5 downloading, copying, storing, and using the Infringed Works to create
6 its training datasets;

7 b. Whether Defendant knew of and materially contributed to the
8 downloading, storage, retention, and copying of the Infringed Works by
9 third parties;

10 c. Whether Plaintiff and members of the Class are entitled to actual and/or
11 statutory damages for the aforementioned violations; and

12 d. Whether Defendant caused further infringement of the Infringed Works
13 by distributing the pre-training dataset under an open license.

FIRST CAUSE OF ACTION
Direct Copyright Infringement,
(17 U.S.C. § 501)
(On Behalf of Plaintiff and the Class)

18 37. Plaintiff repeats the allegations contained in the foregoing paragraphs as
19 if fully set forth herein.

20 38. Plaintiff, as the owner of the registered copyrights of the Infringed
21 Works, holds the exclusive rights to those books under 17 U.S.C. § 106.

22 39. In order to create a comprehensive training dataset known as RedPajama,
23 Together AI downloaded and copied Books3. Books3 includes the Plaintiff's and the
24 members of the Class's Infringed Works. Together AI copied, downloaded, used,
25 stored, and republished many copies of the Books3 dataset and thus Plaintiff's and
26 the Class members' Infringed Works.

27 ||| 40. Neither Plaintiff nor Class members authorized Together AI to make

1 copies, publicly display copies, or distribute copies of their Infringed Works. The U.S.
 2 Copyright Act bestows all the aforementioned rights only on the Plaintiff and the
 3 Class members.

4 41. By downloading, copying, storing, processing, reproducing, and using
 5 the datasets, like Books3, containing copies of Plaintiff's Infringed Works, Defendant
 6 has directly infringed upon Plaintiff's exclusive rights in his copyrighted works.

7 42. By downloading, copying, storing, processing, reproducing, and using
 8 the RedPajama dataset containing copies of Plaintiff's Infringed Works, Defendant
 9 has directly infringed on Plaintiff's exclusive rights in his copyrighted works.

10 43. Defendant repeatedly copied, stored, and used the Infringed Works
 11 without Plaintiff's and members of the Class's permission. Defendant made these
 12 copies without Plaintiff's permission and in violation of their exclusive rights under
 13 the Copyright Act.

14 44. By and through the actions alleged above, Defendant has infringed and
 15 will continue to infringe Plaintiff's copyrights.

16 45. Plaintiff has been injured by Defendant's acts of direct copyright
 17 infringement. Plaintiff is entitled to statutory damages, actual damages, restitution of
 18 profits, and all appropriate legal and equitable relief.

19

20 **SECOND CAUSE OF ACTION**
 21 **Contributory Copyright Infringement**
 22 **17 U.S.C. § 501 *et seq.***
 23 **(On behalf of Plaintiff and the Class)**

24 46. Plaintiff repeats the allegations contained in the foregoing paragraphs as
 25 if fully set forth herein.

26 47. Unknown third parties directly infringed on the Plaintiff's and Class
 27 members' rights by downloading, copying, using, storing and retaining the Infringing
 28 Works. Such use is not fair use and is infringing. Each and every one of these
 infringements is facilitated, encouraged, and made possible by Defendant.

1 48. Defendant had actual and/or red-flag knowledge that the pirated library
2 that it downloaded, copied, stored, duplicated and republished through the RedPajama
3 dataset contained copyrighted works.

4 49. Defendant had actual and/or red-flag knowledge that third parties were
5 downloading and copying its RedPajama dataset that it created using the copyrighted
6 works of Plaintiff and the members of the Class.

7 50. Defendant materially contributed to unknown third party's infringement
8 by: 1) creating the RedPajama dataset which contained copyrighted works; 2)
9 providing its RedPajama dataset to the public for unregulated copying and
10 downloading on both its own specific GitHub and Hugging Face pages; and 3)
11 promoting and encouraging unknown third parties to copy and download its
12 RedPajama dataset.

13 51. Defendant's contribution was substantial and necessary to the scope and
14 persistence of the infringement. Plaintiff and the Class members suffered damages,
15 including statutory damages and profits attributable to the infringement.

PRAYER FOR RELIEF

17 WHEREFORE, Plaintiff, individually and on behalf of all others similarly
18 situated, seeks judgment against Defendant, as follows:

19 a. For an order certifying the Class, naming Plaintiff as Class
20 Representative, and naming Plaintiff's attorneys as Class Counsel to
21 represent the Class;

22 b. For an order declaring that Defendant's conduct violates 17 U.S.C. §
23 501;

24 c. An award of statutory and other damages under 17 U.S.C. § 504 for
25 violations of the copyrights of Plaintiff and the Class by Defendant;

26 d. Reasonable attorneys' fees and reimbursement of costs under 17 U.S.C.
27 §505 or otherwise;

- 1 e. A declaration that such infringement is willful;
- 2 f. Destruction or other reasonable disposition of all copies Defendant made
- 3 or used in violation of the exclusive rights of Plaintiff and the Class,
- 4 under 17 U.S.C. § 503(b);
- 5 g. Pre- and post-judgment interest on the damages awards to Plaintiff and
- 6 the Class, and that such interest be awarded at the highest legal rate from
- 7 and after the date this class action complaint is first served on Defendant;
- 8 h. Further relief for Plaintiff and the Class as the Court deems may be
- 9 appropriate.

10

11 **JURY TRIAL DEMAND**

12 Plaintiff demands a trial by jury on all causes of action and issues so triable.

13

14 DATED: November 6, 2025 Respectfully submitted,

15 DARIUS H. JAMES, individually and on behalf of
16 similarly situated individuals

17 By: /s/ Eugene Y. Turin
18 Plaintiff's Attorney

19 Eugene Y. Turin (SB # 342413)
20 David L. Gerbie (*pro hac vice* forthcoming)
21 Jordan R. Frysinger (*pro hac vice* forthcoming)
22 McGuire Law, P.C.
23 1089 Willowcreek Road, Suite 200
24 San Diego, CA 92131
25 Tel: (312) 893-7002 Ex. 3
26 Fax: 312-275-7895
eturin@mcgpc.com
dgerbie@mcgpc.com
jfrysinger@mcgpc.com

27 *Counsel for Plaintiff and the Putative Class*

EXHIBIT A

Registration Number / Date: TX0003372564 / 1992-08-05

Title: *Negrophobia: an urban parable: a novel* / Darius James.

Copyright Claimant:

Darius James

USCO Catalog Link:

https://publicrecords.copyright.gov/detailed-record/voyager_14732620

Registration Number / Date: TX0004193939 / 1996-01-03

Title: *That's blaxploitation!: roots of the baadassss 'tude* / Darius James.

Copyright Claimant: Darius James

USCO Catalog Link:

https://publicrecords.copyright.gov/detailed-record/voyager_15167875