

BURSOR & FISHER, P.A.

L. Timothy Fisher (State Bar No. 191626)
1990 North California Blvd., 9th Floor
Walnut Creek, CA 94596
Telephone: (925) 300-4455
Facsimile: (925) 407-2700
E-mail: ltfisher@bursor.com

Attorneys for Plaintiff

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

ARTHUR KLEINER, individually and on behalf
of all others similarly situated,

Plaintiff,

v.

ADOBE INC.,

Defendant.

Case No.

CLASS ACTION COMPLAINT

JURY TRIAL DEMANDED

1 Plaintiff Arthur Kleiner (“Plaintiff”) brings this action on behalf of himself and all others
2 similarly situated against Defendant Adobe Inc. (“Defendant”). Plaintiff seeks to recover
3 injunctive relief and damages as a result of Defendant’s unlawful conduct. Plaintiff makes the
4 following allegations pursuant to the investigation of his counsel and are based upon information
5 and belief, except as to the allegations specifically pertaining to himself, which are based on
6 personal knowledge.

7 **NATURE OF THE ACTION**

8 1. This case addresses the surreptitious, non-consensual use and collection of authors’
9 books and written works by Defendant in order to train its SlimLM small language models (SLMs),
10 which is AI software optimized for document assistance tasks on mobile devices. This use violates
11 Adobe’s terms of service at the expense of writers who are unknowingly contributing to training
12 models for Defendant’s SlimLM software.

13 2. SlimLM is a series of small language models created by Defendant that are designed
14 and optimized for mobile deployment.¹ SlimLM, according to Defendant’s own research, is
15 pretrained on a dataset known as “SlimPajama-627B.” It is then fine tuned on DocAssist, a
16 specialized dataset constructed by a team at Adobe.²

17 3. SlimLM, as a language model, is trained by ingesting datasets that are chosen by its
18 developers. It is well-known that developing AI tools necessitates the collection and use of millions
19 of lines of text, including oftencopyrighted materials.³

20 4. The dataset used by Adobe to train its SlimLMs is called SlimPajama, which is a
21 copied, cleaned and deduplicated version of the notorious RedPajama dataset. RedPajama contains
22 the Books3 corpus, which contains hundreds of thousands of copyrighted books that were acquired
23 without the authorization or consent of the authors. The SlimPajama dataset also contains content
24

25 ¹ Adobe Research, *SlimLM: An Efficient Small Language Model for On-Device Document Assistance* (August 1, 2025).

26 ² *Id.*

27 ³ Ars Technica, *OpenAI says it’s “impossible” to create useful AI models without copyrighted material* (Jan. 9, 2024) <https://arstechnica.com/information-technology/2024/01/openai-says-its-impossible-to-create-useful-ai-models-without-copyrighted-material/>
28

1 from the Common Crawl dataset, which is also known to contain copyrighted and unauthorized
2 material.⁴

3 5. Plaintiff and Class members are authors and copyright holders. They own registered
4 copyrights in books that were included in the SlimPajama (via Common Crawl and RedPajama)
5 datasets that Defendant downloaded, copied, stored, and used without their permission or
6 compensation.

7 6. Plaintiff and Class members never authorized Defendant to download, copy, store,
8 or use their copyrighted works. Defendant has never compensated Plaintiffs and Class members for
9 downloading, copying, storing, or using their copyrighted works.

10 7. Adobe has and continues to benefit commercially from its massive acts of copyright
11 infringement. It does so by securing lucrative contracts with enterprise customers for the use of its
12 AI software, including through the SlimLM AI platform, and through deploying such AI software
13 (or tools relying on SlimLM) through its suite of Adobe branded products.

14 8. Through the above acts, Defendant has infringed Plaintiff's copyrighted works and
15 continues to do so by continuing to store, copy, use, and process the datasets containing copies of
16 Plaintiff's and the putative Class's copyrighted books.

17 **PARTIES**

18 9. Plaintiff Arthur Kleiner is a citizen of New York and resident of Manhattan, New
19 York. Plaintiff Kleiner is a published author of many books, including *The Age of Heretics:*
20 *Heroes, Outlaws, and the Forerunners of Corporate Change*, which was published in 1996.

21 10. Plaintiff Kleiner registered his book with the United States Copyright Office in
22 1995, and has held the ownership of the copyright since.

23 11. Plaintiff Kleiner's book, *The Age of Heretics*, was included in the datasets that
24 Defendant pirated, copied used, and transcribed to train its SlimLM models.

25 _____
26 ⁴ Rights Alliance for the Creative Industries on the Internet, *Report on AI model providers' training*
27 *data transparency and enforcement of copyrights*, [https://rettighedsalliancen.com/wp-](https://rettighedsalliancen.com/wp-content/uploads/2024/09/Report-on-AI-model-providers-training-data-transparency-and-enforcement-of-copyrights.pdf)
28 [content/uploads/2024/09/Report-on-AI-model-providers-training-data-transparency-and-](https://rettighedsalliancen.com/wp-content/uploads/2024/09/Report-on-AI-model-providers-training-data-transparency-and-enforcement-of-copyrights.pdf)
[enforcement-of-copyrights.pdf](https://rettighedsalliancen.com/wp-content/uploads/2024/09/Report-on-AI-model-providers-training-data-transparency-and-enforcement-of-copyrights.pdf) (Sep. 5, 2024); *see also* The Atlantic, *The Company Quietly*
Funneling Paywalled Articles to AI Developers (Nov. 4, 2025)
<https://www.theatlantic.com/technology/2025/11/common-crawl-ai-training-data/684567/>

1 18. The Pile is a dataset curated by a research organization called EleutherAI for use in
2 training AI models. In December 2020, EleutherAI introduced this dataset in a paper called “The
3 Pile: An 800GB Dataset of Diverse Text for Language Modeling.”⁵ The paper provides a
4 description of the Books3 dataset contained within The Pile:

5 Books3 is a dataset of books derived from a copy of the contents of the Bibliotik
6 private tracker ... Bibliotik consists of a mix of fiction and nonfiction books and
7 is almost an order of magnitude larger than our next largest book dataset
8 (BookCorpus2). We included Bibliotik because books are invaluable for long-
range context modeling research and coherent storytelling.⁶

9 19. Bibliotik is one of a number of notorious pirate websites that also includes Library
10 Genesis (aka LibGen), Z-Library (aka B-ok), Sci-Hub, and Anna’s Archive. These illegally
11 sourced libraries have long been of interest to the AI-training community because they contain vast
12 quantities of unauthorized copyrighted material, including books, which are required to train
13 LLM’s (large language models) and SLMs (small language models). Using such illegal libraries
14 permits AI companies like Defendant to gain quick access to large quantities of training data
15 rapidly. For that reason, these pirate libraries also violate the U.S. Copyright Act.

16 20. The person who assembled the Books3 dataset, Shawn Presser, has confirmed in
17 public statements that Books3 represents “all of Bibliotik” and contains approximately 196,640
18 books.

19 21. EleutherAI’s website encourages the public to download The Pile dataset from a
20 website known as “The-Eye.” Anyone who downloads The Pile from The-Eye is therefore also
21 downloading a copy of Books3.

22 22. The Pile is also available for download from the “Hugging Face” website. Before
23 October 2023, the Books3 subset of The Pile was available for download from Hugging Face as a
24 standalone dataset. But in October 2023, the Books3 dataset was removed with a message that it “is
25 defunct and no longer accessible due to reported copyright infringement.”⁷

26 ⁵ <https://arxiv.org/pdf/2101.00027.pdf>

27 ⁶ *Id.* at 3–4.

28 ⁷ https://web.archive.org/web/20231127101818/https://huggingface.co/datasets/the_pile_books3

1 23. Books3 has also been included in its entirety within another dataset known as
2 RedPajama, which was created by the company Together AI. Released in or around April 2023, the
3 RedPajama dataset contained a subset called “Books” or “RedPajama-Books” that was a direct
4 copy of the “Books3 dataset.” The RedPajama dataset is available for download from Hugging
5 Face. Before October 2023, the Books3 subset was “downloaded from Huggingface [sic]” when a
6 user ran the script that automatically assembled the RedPajama dataset.⁸ After the “Books3
7 dataset” was removed from Hugging Face in October 2023, the RedPajama dataset documentation
8 similarly added a message that Books3 is defunct “due to reported copyright infringement.”⁹

9 24. But before October 2023, anyone who downloaded the RedPajama or The Pile
10 datasets from Hugging Face was necessarily downloading a copy of the Books3 dataset.

11 25. SlimPajama is a “deduplicated, multi-corpora” dataset that was created by “cleaning
12 and deduplicating the 1.21T token RedPajama dataset.”¹⁰

13 26. In this context, “cleaning” data does not mean removing copyrighted or protected
14 works. It means removing errors, corrupted data, or duplicated/repeated data, with the result of
15 making a dataset more accurate and streamlined.

16 27. This means that the SlimPajama dataset contains copyrighted material that has not
17 been licensed and authorized by its creators.

18
19
20
21
22
23 _____
24 ⁸ <https://web.archive.org/web/20230420075601/https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T>

25 ⁹ <https://web.archive.org/web/20240510231649/https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T>

26 ¹⁰ Cerebras.ai, *SlimPajama: A 627B token, cleaned and deduplicated version of RedPajama*,
27 <https://www.cerebras.ai/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama> (Jun. 9, 2023).
28

1 28. The SlimPajama dataset is also based in part on—or overlaps with—the Common
2 Crawl dataset,¹¹ which is known to contain copyrighted material, including “articles published on
3 news websites and in digital newspapers.”¹²

4 29. Plaintiff’s copyrighted book, listed in **Exhibit A**, is among the works in the Books3
5 dataset—and therefore the RedPajama, and SlimPajama— dataset. Below, this book is referred to
6 as the **Infringed Work**.

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

8 30. A large language model (“LLM”) is Artificial Intelligence software designed to emit
9 convincingly naturalistic text outputs in response to user prompts. Small language models (“SLM”) are similar to LLMs, but trained “on data from specific domains,” meaning their training utilizes
10 smaller or more targeted sets of data.¹³

11 31. Though SLMs and LLMs are software programs, they are not created the way most
12 software programs are—that is, by human software programmers writing code. Rather, an SLM is
13 *trained by* copying extremely large quantities of textual works and then feeding these copies into
14 the model. This corpus of input material is called the *training dataset*.

15 32. Training consists of a multi-stage process (known as the training pipeline) that
16 includes the acquisition and curation of the dataset, processing of the dataset, feeding the dataset
17 into the model so that the model can extract the patterns and relationships from the protected
18 expression contained therein; and further fine-tuning the model for more specialized uses with even
19 more data. This process also involves experimentation to improve the data mixture (i.e., the final
20

21 _____
22 ¹¹ Rights Alliance for the Creative Industries on the Internet, *Report on AI model providers’*
23 *training data transparency and enforcement of copyrights*, [https://rettighedsalliancen.com/wp-](https://rettighedsalliancen.com/wp-content/uploads/2024/09/Report-on-AI-model-providers-training-data-transparency-and-enforcement-of-copyrights.pdf)
24 *content/uploads/2024/09/Report-on-AI-model-providers-training-data-transparency-and-*
25 *enforcement-of-copyrights.pdf* (Sep. 5, 2024); *see also* The Atlantic, *The Company Quietly*
Funneling Paywalled Articles to AI Developers (Nov. 4, 2025); *see also* SlimPajama-DC:
Understanding Data Combinations for LLM Training,
<https://ar5iv.labs.arxiv.org/html/2309.10818#S3.SS1>

26 ¹² *Netherlands: Common Crawl removes 2 million scraped news articles, complying with request*,
27 *Piracy Mon*[https://piracymonitor.org/netherlands-common-crawl-removes-2-million-scraped-](https://piracymonitor.org/netherlands-common-crawl-removes-2-million-scraped-news-articles-complying-with-request/)
news-articles-complying-with-request/

28 ¹³ *LLMs vs. SLMs: The Differences in Large & Small Language Models*, Splunk,
https://www.splunk.com/en_us/blog/learn/language-models-slm-vs-llm.html (Feb. 17, 2025).

1 training dataset and the proportion of each component dataset) of the model. Each experiment is
2 also known as an “ablation” in technical parlance.

3 33. The first step in training the model is acquiring and curating the data that goes in to
4 the model. This acquisition necessarily includes copying and downloading data, usually onto
5 persistent storage. This is because training an SLM or any other generative AI model is expensive
6 and resource intensive, so retaining persistent copies of datasets avoids the need to reacquire the
7 data when developing new models, or altering the data mixtures during the development process.

8 34. Training an SLM is not only a function of quantity of data, but also of quality. The
9 selection and curation of training data is therefore an important first step in training. Copyrighted
10 books are well known among developers of AI models to be high-quality data for training SLMs.

11 35. After the initial copying and processing of data, the SLM copies each textual work
12 in the training dataset and extracts protected expression from it. During what is known as
13 *pretraining*, the SLM progressively adjusts its output to more closely approximate the protected
14 expression copied from the training dataset. The SLM records the results of this process in a large
15 set of numbers called weights (also known as *parameters*) that are stored within the model. These
16 weights are entirely and uniquely derived from the protected expression in the training dataset.
17 Once a model is pretrained, the result is a trained model known as a *base* or *foundational* model.

18 36. During the development process of an SLM, engineers may also conduct
19 experiments known as ablations or “ablation studies” that test the effect of certain data on the
20 model. This can include, for example, determining whether there is a difference in the quality of a
21 model’s output if it is trained with certain books or without. A dataset may be used to run such
22 experiments but ultimately be excluded from the final training mixture of the model. Importantly,
23 these datasets used for ablation studies often also consist of copyrighted works, including books.

24 37. Once an SLM has copied the textual works in the training dataset and extracted the
25 protected expression into stored weights, an SLM is able to emit convincing simulations of natural
26 written language in response to user prompts. Whenever an SLM generates text output in response
27 to a user prompt, it is performing a computation that relies on these stored weights which represent
28 the protected expression ingested from the training dataset.

1 38. Throughout each step of the training pipeline, the same dataset is typically used
2 (i.e., copied) multiple times. Indeed, given the cost of developing an SLM, it is a ubiquitous
3 practice to retain datasets for future use, whether that use is to pretrain future models, to perform
4 ablations on a model, or to fine-tune an already trained base model. Each step involves making
5 additional copies of the underlying data. The implication is that if a dataset contains unlawfully
6 obtained copyrighted material, each step of the training pipeline may result in an unauthorized use
7 (i.e., infringement) of that copyrighted work.

8 **C. Adobe Used Datasets Known to Contain Pirated Books to Develop its Small Language**
9 **Models**

10 39. Adobe is a computer software company that “empowers everyone, everywhere to
11 imagine, create, and bring any digital experience to life. From creators and students to small
12 businesses, global enterprises, and nonprofit organizations — customers choose Adobe products to
13 ideate, collaborate, be more productive, drive business growth, and build remarkable
14 experiences.”¹⁴

15 40. Adobe’s software includes the popular PDF editing suite known as Adobe Acrobat,
16 the image editor software known as Photoshop, the video editor software known as Adobe
17 Premiere, and the graphic design tool known as Adobe Illustrator.

18 41. As part of Adobe’s improvement of its “consumer technology” products, Defendant
19 has developed SlimLM, “a series of small language models specifically designed and optimized for
20 mobile deployment.”¹⁵

21 42. Adobe says that its SlimLM SLM is pretrained on SlimPajama-627B. As detailed
22 above, SlimPajama is known to contain the contents of RedPajama, which in turns contains part of
23 the Books3 corpus, comprised of thousands of copyrighted works.

24
25
26 ¹⁴ Adobe, *Changing the world through personalized digital experiences*,
27 <https://www.adobe.com/about-adobe.html> (last accessed Jan. 22, 2026).

28 ¹⁵ Adobe Research, *SlimLM: An Efficient Small Language Model for On-Device Document Assistance* (August 1, 2025).

1 43. SlimPajama is also known to contain parts of, if not all, of the Common Crawl
2 dataset. Specifically, SlimPajama’s composition (by token percentage) is 52.2% Common Crawl.¹⁶

3 44. The Common Crawl dataset was created and curated by Common Crawl. “Common
4 Crawl is a non-profit organization that crawls the internet to provide an extensive archive or dataset
5 to the public. Tech companies use this dataset as training data for their generative AI models,
6 including Apple’s openELM, Microsoft’s Phi, OpenAI’s ChatGPT, NVIDIA’s Nemo Megatron,
7 Deepseek’s Deepseek V3, and Anthropic’s Claude.”¹⁷

8 45. BREIN, a Dutch non-profit foundation, asked the creators of Common Crawl to
9 remove over “two million news articles belonging to popular Dutch news outlets from its AI
10 training dataset.” According to BREIN, “these articles were copied without permission.”¹⁸

11 46. Common Crawl is known to contain other copyrighted works beyond the millions of
12 unlicensed Dutch news articles.¹⁹

13 47. In other words, Defendant’s SlimLM is trained on potentially multiple datasets
14 known to contain copyrighted works: SlimPajama and Common Crawl.

15 48. This runs contrary to Adobe’s outspoken stance on ethical AI development and use.

16 49. In numerous places on its website, Adobe touts its Adobe Firefly product, an AI
17 image generator that was built through “thoughtful, responsible development.” Adobe claims it
18 “developed Adobe Firefly to prevent it from creating content that infringes copyright or intellectual
19 property rights, and it is designed to be commercially safe,” and that they “do not mine content
20 from the web to train Adobe Firefly.”²⁰

21
22
23 ¹⁶ Emergent Mind, *SlimPajama Dataset*, <https://www.emergentmind.com/topics/slimpajama-dataset>

24 ¹⁷ CyberNews, *Common Crawl removes AI dataset containing over 2M news articles*,
25 <https://cybernews.com/ai-news/common-crawl-ai-dataset-2m-news-articles/> (Nov. 5, 2025).

26 ¹⁸ *Id.*

27 ¹⁹ The Business Standard, *New York Times successfully removes copyrighted content from AI training dataset*, <https://www.tbsnews.net/tech/new-york-times-successfully-removes-copyrighted-content-ai-training-dataset-741266> (Nov. 18, 2023).

28 ²⁰ <https://www.adobe.com/ai/overview/ethics.html>

1 50. Despite Adobe’s claims of ethical AI development for its Firefly product, Adobe
2 has nevertheless celebrated the development of its SlimLM SLM, which was trained on multiple
3 datasets of copyrighted material, including Plaintiff’s. At no point in Adobe’s paper describing its
4 SlimLM model did it discuss compensation for the writers and creators of the works on which
5 SlimLM was trained.²¹

6 51. Adobe infringed on Plaintiff and Class members’ copyrighted works on a massive
7 scale. Adobe downloaded these books from the SlimPajama dataset—which includes copyrighted
8 works via RedPajama and Common Crawl—without authorization from or compensation to their
9 authors. Adobe then continued copying and storing the datasets, and used them to develop and train
10 the SlimLM series of SLMs.

11 52. Adobe stored and used Plaintiff’s and Class members’ copyrighted works to train
12 other models, including non-public models, whether sourced from SlimPajama, Common Crawl or
13 other pirated and/or unauthorized, copyrighted sources.

14 53. Plaintiff and Class members have retained ownership rights in their works. Plaintiff
15 and Class members did not consent to the use of their works as training material for SlimLM.
16 Nonetheless, their works were downloaded, copied, stored, and used for training SlimLM to be a
17 critical support system for current or future Adobe products.

18 **D. Adobe’s Use of Copyrighted Works Is Not Fair Use.**

19 54. Adobe’s downloading, copying, storage, and use of Plaintiff’s and Class members’
20 works was not fair use. Adobe could have—but chose not to—lawfully obtain the Infringed Work.
21 It is near impossible that “any accused infringer could ever meet its burden of explaining why
22 downloading source copies from pirate sites *that it could have purchased or otherwise accessed*
23 *lawfully* was itself reasonably necessary to any subsequent fair use.” *Bartz v. Anthropic PBC*, 787
24 F. Supp. 3d 1007, 1025 (N.D. Cal. 2025). Adobe used the SlimPajama dataset (comprised of data
25 from the RedPajama and Common Crawl datasets), which includes materials sourced from shadow
26 libraries known to indiscriminately Hoover up countless copyrighted works, for its centralized

27 _____
28 ²¹ Adobe Research, *SlimLM: An Efficient Small Language Model for On-Device Document Assistance* (August 1, 2025)

1 commercial database, all with complete disregard to the rights of copyright holders. Adobe used
2 these datasets containing Plaintiff's and Class members' copyrighted works as a substitute for
3 purchasing or compensating creators for authorized copies of those works. In doing so, Adobe
4 displaced or diluted the market for Plaintiffs' and Class members' works. Plaintiff's works were
5 copied and maintained for future commercial purposes, including to further develop its SlimLM
6 series of SLMs. Pirating otherwise purchasable works is copyright infringement.

7 55. The exploitation of Plaintiff's Infringed Work was not indirect. It was and is a direct
8 and illegal download, use and scraping of the entirety of Plaintiff's Infringed Work (as part of the
9 SlimPajama dataset) with no transformation of form. Plaintiff's Infringed Work was copied, used
10 and maintained for future commercial purposes. Defendant's use was not transitory as Plaintiff's
11 Infringed Work was not immediately destroyed, but was retained to train Adobe's SLMs.
12 Integration of this technology into Adobe's existing document editing and creation suite could
13 generate billions of dollars for Adobe. Adobe's "piracy of otherwise available copies is inherently,
14 irredeemably infringing even if the pirated copies are immediately used for [a] transformative use
15 and immediately discarded." *Id.* at 1025.

16 56. Adobe's use of these copyrighted works to train its SlimLM SLMs are not
17 transformative. SLM outputs are the result of statistical outputs based on the training corpora.
18 Broadly, SLMs operate by probabilistically predicting the next "token." Tokens are units of
19 language which can be words, combinations of letters, or even a punctuation or space. Adobe's
20 SlimLM models are no different. Recent studies demonstrate that the LLMs created by the leading
21 AI developers are capable of regurgitating substantial portions of the training data for those
22 models, even with the implementation of safeguards designed to prevent the regurgitation of
23 training data.²²

24 57. Even if the court applies the statutory fair use factors to Adobe's unmitigated
25 copying through downloading datasets including SlimPajama, and its use of those datasets to train

26 _____
27 ²² Ahmed Ahmed et. a., "Extracting Books from Production Language Models" (Jan. 6, 2026),
28 available at <https://arxiv.org/pdf/2601.02671>. This study analyzed the output of LLMs created by Anthropic, Google, OpenAI, and xAI—models which use the same or similar architecture to the xGen models.

1 its SLMs, the application of the factors weighs against fair use. The four factors of fair use are: (1)
2 “the purpose and character of the use,” (2) “the nature of the copyrighted work,” (3) “the amount
3 and substantiality of the portion used,” and (4) “the effect of the use upon [Plaintiff’s] potential
4 market.” 17 U.S.C. § 107.

5 58. **Factor one.** The purpose and character of Adobe’s use is strictly commercial and
6 not transformative. The factor one analysis includes (1) whether the purpose of copying is
7 commercial or non-commercial and (2) whether the purpose is transformative. Adobe’s use was
8 and is not for educational or commentary purposes, but furthered a commercial endeavor to utilize
9 a massive library of valuable copyrighted works that helped train and develop financially lucrative
10 Adobe products. Further, there is nothing transformative about using unauthorized copies of
11 Plaintiff’s Infringed Work to train Adobe’s SlimLM models in place of sourcing free market
12 alternatives. Piracy is not transformative and is not fair use.

13 59. **Factor two.** The factor two analysis primarily questions whether the nature of the
14 copied work is more creative or factual, and secondarily questions whether the copied works are
15 accessible or out of print. Here, Plaintiff’s work is highly creative and is therefore entitled to
16 broader protections under copyright. Further, inaccessible works are provided less protection
17 because there is no market for inaccessible works that a defendant’s copying could displace.
18 However, “this is not a case where source copies were unavailable,” *Anthropic PBC*, 787 F. Supp.
19 3d at 1027, and Adobe chose not to purchase or license this copyrighted material – even though
20 that option was readily available. The nature of Plaintiff’s Infringed Work ensures Plaintiff the
21 greatest copyright protection and narrows Defendant’s argument for fair use.

22 60. **Factor three.** The third factor asks whether a qualitatively and/or quantitatively
23 substantial portion of Plaintiff’s Infringed Work was stolen. The entirety of Plaintiff’s Infringed
24 Work was compiled in the SlimPajama dataset, and copied and stored for Adobe’s use on an
25 ongoing basis. These works were not copied in a modified form (for example, as an image
26 thumbnail) into the dataset, such that the substance was not extracted. Plaintiff’s Infringed Work
27 was copied in its entirety to take both the whole and the heart of the work. The third factor weighs
28 strongly in favor of infringement.

1 65. Specifically excluded from the Class is Defendant, Defendant’s officers, directors,
2 agents, trustees, parents, children, corporations, trusts, representatives, employees, principals,
3 servants, partners, joint ventures, or entities controlled by Defendant, and its heirs, successors,
4 assigns, or other persons or entities related to or affiliated with Defendant and/or Defendant’s
5 officers and/or directors, the judge assigned to this action, and any member of the judge’s
6 immediate family.

7 66. Plaintiff reserves the right to expand, limit, modify, or amend the class definition,
8 including the addition of one or more subclasses, in connection with his motion for class
9 certification, or at any other time, based on, *inter alia*, changing circumstances and/or new facts
10 obtained.

11 67. **Numerosity.** On information and belief, hundreds of thousands of copyright
12 holders fall into the definition of the Class. Members of the Class can be identified through
13 Defendant’s records, discovery, and third-party sources.

14 68. **Commonality and Predominance.** Common questions of law and fact exist as to
15 all members of the Class and predominate over any questions affecting only individual members of
16 the Class. These common legal and factual questions include, but are not limited to, the following:

- 17 (a) Whether Defendant violated the copyrights of Plaintiff and the Class when it
18 obtained copies of Plaintiff’s Infringed Work, stored the Infringed Works,
19 copied the Infringed Works, and used the Infringed Works to develop
20 Defendant’s SLMs;
- 21 (b) Whether any affirmative defense excuses Defendant’s conduct; and
- 22 (c) Whether any statutes of limitation limit the potential for recovery for
23 Plaintiffs and the Class.

24 69. **Typicality.** Plaintiff’s claims are typical of the claims of the other members of the
25 Class in that, among other things, all Class members were similarly situated and were comparably
26 injured through Defendant’s wrongful conduct as set forth herein. Further, there are no defenses
27 available to Defendant that are unique to Plaintiff.

1 75. Adobe downloaded, ingested, or otherwise acquired copies of the SlimPajama
2 dataset containing, in part, Common Crawl and Redpajama, which includes the Infringed Work.
3 The initial downloading constitutes the first unauthorized copying of Plaintiff’s works by Adobe in
4 the training pipeline.

5 76. Upon information and belief, Adobe created and stored numerous copies of
6 SlimPajama in its internal servers.

7 77. Upon information and belief, to develop the SlimLM SLMs, Adobe copied the
8 SlimPajama dataset containing Common Crawl and RedPajama (comprised in part by Books3) to
9 develop these models and incorporate the dataset (and various iterations thereof) into the models’
10 training data.

11 78. Plaintiff and the Class members never authorized Defendant to make copies of their
12 Infringed Works, make derivative works, publicly display copies (or derivative works), or
13 distribute copies (or derivative works). All those rights belong exclusively to Plaintiff and the Class
14 members under the U.S. Copyright Act.

15 79. By copying, storing, processing, reproducing, and using the dataset containing
16 pirated copies of Plaintiff’s Infringed Works, Defendant has directly infringed Plaintiff’s exclusive
17 rights in their copyrighted works.

18 80. By copying, storing, processing, and reproducing the SlimLM models trained on
19 Plaintiff’s Infringed Work, Adobe has directly infringed Plaintiff’s exclusive rights in their
20 copyrighted works.

21 81. Defendant repeatedly copied, stored, and used the Infringed Work without
22 Plaintiff’s permission. Defendant made these copies without Plaintiff’s permission and in violation
23 of their exclusive rights under the Copyright Act.

24 82. Because the complete training mixture for the SlimLM models has not been made
25 public, other unauthorized repositories of pirated copies of Plaintiff’s and Class members’
26 copyrighted works may have been used to develop Adobe’s SlimLM and other models.
27
28

- (g) Pre- and post-judgment interest on the damages awarded to Plaintiff and the Class, and that such interest be awarded at the highest legal rate from and after the date this Complaint is first served on Defendant.
- (h) For an order establishing that Defendant is 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) For injunctive relief as the Court may deem proper; and
- (j) Further relief for Plaintiff and the Class as may be appropriate.

JURY TRIAL DEMANDED

Plaintiff demands a trial by jury on all claims so triable.

Dated: February 9, 2026

BURSOR & FISHER, P.A.

By: /s/ L. Timothy Fisher

L. Timothy Fisher (State Bar No. 191626)
1990 North California Blvd., 9th Floor
Walnut Creek, CA 94596
Telephone: (925) 300-4455
Facsimile: (925) 407-2700
E-mail: ltfisher@bursor.com

Attorneys for Plaintiff

EXHIBIT A

Copyright

Registration Number / Date:

TX0004280899 / 1996-05-31

Type of Work:

Text

Title:

The age of heretics :heroes, outlaws, and the forerunners of corporate change /Art Kleiner.

Date of Creation:

1995

Date of Publication:

1996-05-01

Copyright Claimant:

Art Kleiner

Description:

414 p.

Imprint:

New York : Currency, c1996.

Edition:

1st ed.

Names:

Kleiner, Art

USCO Catalog Link:

https://publicrecords.copyright.gov/detailed-record/voyager_15216775

Disclaimer: This material was generated by the U.S. Copyright Office's Copyright Public Records System (CPRS). For certified records, contact the [Records Research and Certification Division](#). For information on searching copyright records, see [How to Investigate the Copyright Status of a Work \(Circular 22\)](#). For information on removing personal information from Copyright Office public records, refer to [Privacy: Public Copyright Registration Records\(Circular 18\)](#).