

1 Joseph R. Saveri (SBN 130064)
Diane S. Rice (SBN 118303)
2 Christopher K.L. Young (SBN 318371)
William W. Castillo Guardado (SBN 294159)

3 **SAVERI LAW FIRM, LLP**
550 California Street, Suite 910
4 San Francisco, CA 94104
Telephone: (415) 500-6800
5 Facsimile: (415) 395-9940
Email: jsaveri@saverilawfirm.com
6 drice@saverilawfirm.com
7 cyoung@saverilawfirm.com
wcastillo@saverilawfirm.com

8 *Attorneys for Individual and Representative Plaintiff*
9 *E. Molly Tanzer*

10 **UNITED STATES DISTRICT COURT**
11 **NORTHERN DISTRICT OF CALIFORNIA**
12 **SAN FRANCISCO DIVISION**

14 E. MOLLY TANZER, an individual,
15 Individual and Representative Plaintiff,
16 v.
17 ADOBE INC.
18 Defendant.

Case No.
COMPLAINT
JURY TRIAL DEMANDED
CLASS ACTION

TABLE OF CONTENTS

1

2 OVERVIEW1

3 JURISDICTION AND VENUE2

4 PARTIES.....2

5 A. PLAINTIFF2

6 B. DEFENDANT3

7 C. AGENTS AND CO-CONSPIRATORS3

8 FACTUAL ALLEGATIONS3

9 CLASS ALLEGATIONS17

10 CAUSES OF ACTION19

11 COUNT I.....19

12 DEMAND FOR JUDGMENT.....21

13 JURY TRIAL DEMANDED22

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

1 Plaintiff E. Molly Tanzer (“Plaintiff”), on behalf of herself and all others similarly situated,
2 brings this Complaint against Defendant Adobe Inc. (“Adobe”).

3 **OVERVIEW**

4 1. Adobe engaged in a massive piracy scheme to develop the Nemotron and SlimLM
5 family of large language models (“LLMs”). Specifically, Adobe partnered with NVIDIA to develop
6 “next-generation NVIDIA AI foundation models,” including the Nemotron family of models trained
7 using millions of books. No public information exists showing that Adobe or NVIDIA licensed
8 millions of books to train their AI models. Upon information and belief, Adobe and NVIDIA were
9 only able to obtain such a large number of books by sourcing them from pirated sources, such as
10 Anna’s Archive, a shadow library that purposefully flouts copyright law and illegally distributes
11 millions of copyrighted books.

12 2. Adobe also acquired the notorious Books3 dataset of nearly 200,000 pirated books to
13 develop its SlimLM LLMs. The Books3 dataset is a subset of the SlimPajama dataset that Adobe
14 acquired. Books3 contains the entirety of copyrighted books from the Bibliotik shadow library that
15 were obtained without the authorization or consent of the authors.

16 3. Plaintiff and Class members are authors. They own registered copyrights in certain
17 books that were included in datasets originating from Books3 and Anna’s Archive that were used by
18 Defendant and which were copied, stored, and used without their permission or compensation.

19 4. Adobe admits that it must compensate rightsholders, obtain their authorization, and
20 respect their copyrights before including their works in training data for AI models. For example,
21 Adobe’s Firefly AI models are trained on *licensed* images whose creators are paid on a yearly basis for
22 the inclusion of those images in Firefly’s training data. In contrast, book authors have *never* received
23 compensation for Adobe’s use, storage, or copying of their works, nor have these authors ever
24 provided Adobe with authorization to use, copy, or store their works. Adobe’s claims of “respecting”
25 the creative community ring hollow.

26 5. Adobe benefitted commercially from its acts of massive copyright infringement,
27 including by securing contracts with enterprise customers for use of its LLMs, and by incorporating
28

1 the infringing Nemotron and SlimLM models into Adobe programs and tools, including Adobe
2 Acrobat.

3 6. Through the above acts, Defendant has infringed Plaintiff’s copyrighted works, and it
4 continues to do so by continuing to store, copy, use, and process datasets containing copies of
5 Plaintiff’s and the putative Class’s copyrighted books.

6 **JURISDICTION AND VENUE**

7 7. This Court has subject-matter jurisdiction under 28 U.S.C. § 1331 and §1338(a) because
8 this is a civil action arising under the copyright laws of the United States, including 17 U.S.C. § 501.

9 8. Jurisdiction and venue are proper in this judicial district under 28 U.S.C. § 1391(c)(2)
10 because Adobe is headquartered in this district. Adobe acquired the Books3 dataset to create the
11 SlimLM series of large language models, and, upon information and belief, developed the Nemotron
12 models using datasets sourced from the shadow library Anna’s Archive. Defendant commercially
13 benefited from the creation of these models. Therefore, a substantial part of the events giving rise to
14 the claim occurred in this District. A substantial portion of the affected interstate trade and commerce
15 was carried out in this District. Defendant has transacted business, maintained substantial contacts,
16 and/or committed overt acts in furtherance of the illegal scheme and conspiracy throughout the United
17 States, including in this District. Defendant’s conduct has had the intended and foreseeable effect of
18 causing injury to persons residing in, located in, or doing business throughout the United States,
19 including in this District.

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

23 **PARTIES**

24 **A. PLAINTIFF**

25 10. Plaintiff E. Molly Tanzer is an author residing in Colorado who owns registered
26 copyrights in multiple books, including *Creatures of Will and Temper*, *Creatures of Charm and*
27 *Hunger*, and *Vermilion*.

1 11. A non-exhaustive list of registered copyrights owned by Plaintiff is included as Exhibit

2 A.

3 **B. DEFENDANT**

4 12. Defendant Adobe Inc. is a Delaware corporation headquartered at 345 Park Avenue,
5 San Jose, CA 95110.

6 **C. AGENTS AND CO-CONSPIRATORS**

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

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

17 **FACTUAL ALLEGATIONS**

18 15. Generative artificial intelligence (“AI”) models are designed to generate output—such as
19 text, images, video, or audio—in response to user prompts. An AI model is not created the way most
20 software programs are—that is, by human software programmers writing code. Rather, an AI model is
21 *trained* by copying an enormous quantity of data and then feeding these copies into the model. This
22 corpus of input material is called the *training dataset*.

23 16. These models are trained with vast quantities of data to allow the output to closely
24 recreate the content of its training data. For example, an image model instructed to generate a picture of
25 a cat will emit a picture based on pictures of cats used to train the model.

1 17. A common type of generative AI model is known as a large language model (“LLM”)
2 designed to emit text outputs in response to user prompts. Often, the models are created to generate text
3 mimicking human language.

4 18. Training consists of a multi-stage process that includes the acquisition and curation of
5 the dataset, processing of the dataset, feeding the dataset into the model so that the model can extract
6 the patterns and relationships from the protected expression contained therein, and further fine-tuning
7 the model for more specialized uses with even more data.

8 19. The first step in training the model is acquiring and curating the data that goes into the
9 model. Training an AI model is not only a function of quantity of data, but also of quality. As
10 Adobe’s own General Counsel Dana Rao acknowledged, generative AI “is only as good as the data on
11 which it’s trained.”¹ The selection and curation of training data is therefore an important first step in
12 training.

13 20. Copyrighted books tend to be high-quality data for training LLMs. Studies have found
14 models to perform best when their training data comprises books,² such that when books are removed
15 from the training mixture of AI models, model performance degrades more than with the removal of
16 other types of training data.³ Put differently, the inclusion of books in the training data mixture has a
17 disproportionately large impact on a model’s capabilities relative to their share of training tokens.⁴ For
18 these reasons, the inclusion of books in the training mixture of LLMs is *essential*.

19
20
21
22 ¹ Rachel Metz & Brody Ford, *Adobe’s ‘Ethical’ Firefly AI Was Trained on Midjourney Images*, YAHOO
23 FINANCE (Apr. 12, 2024), <https://finance.yahoo.com/news/adobe-ethical-firefly-ai-trained-123004288.html> (on file with counsel).

24 ² Jack W. Rae, et al., *Scaling Language Models: Methods, Analysis & Insights from Training Gopher*,
(Jan. 21, 2022) [<https://arxiv.org/abs/2112.11446>].

25 ³ Shayne Longpre, et al., *A Pretrainer’s Guide to Training Data: Measuring the Effects of Data Age, Domain Coverage, Quality, & Toxicity*, in Proceedings of the 2024 Conference of the North American
26 Chapter of the Association for Computational Linguistics (NAACL) 3245–3276, 3246 (2024).

27 ⁴ Yang Zhao et al., *Deciphering the Impact of Pretraining Data on Large Language Models through Machine Unlearning*, in *Findings of the Association for Computational Linguistics: ACL 2024* 9386–
28 9406 (2024).

1 21. During training, the AI model copies and ingests each work in the training dataset and
2 extracts protected expression from it. During what is known as *pretraining*, the AI model progressively
3 adjusts its output to more closely approximate the protected expression copied from the training
4 dataset. The AI model records the results of this process in a large set of numbers called weights (also
5 known as *parameters*) that are stored within the model. These weights are entirely and uniquely derived
6 from the protected expression in the training dataset. Once a model is pretrained, it results in a trained
7 model known as a *base* or *foundational* model.

8 22. After a model is pretrained, it may then be *fine-tuned* with additional data. According to
9 Adobe:

10 During the build and train phase of AI development, training is not a “one
11 and done” process. The model typically must be “tuned,” that is, trained on
12 either a larger, more diverse dataset or on a completely different dataset than
13 originally expected. For example in the initial training of Adobe Firefly, the
14 product team realizing that they needed more images of hands in the dataset
15 to improve the model.⁵

16 23. Engineers may also conduct experiments known as “ablation studies” that test the effect
17 of certain data on the model. This can include, for example, determining whether there is a difference in
18 the quality of a model’s output if it is trained with books versus without. A dataset may be used to run
19 such experiments but be excluded from the final training dataset of the model. Importantly, these
20 datasets too may consist of copyrighted works, including books.

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

26 25. Recent studies confirm that LLMs are capable of *memorizing* their training data. LLMs
27 are able to emit verbatim or reproduce near-verbatim portions of their training data. One study found

28 ⁵ *Adobe Generative AI Built for Business Solution Brief*, ADOBE (2024), <https://www.adobe.com/cc-shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-brief.pdf> (on file with counsel).

1 that the LLMs of several leading AI developers—Anthropic, Google, OpenAI, and xAI—can
2 regurgitate substantial portions of their training data, even with the implementation of safeguards
3 designed to prevent the regurgitation of training data.⁶ Another recent study found that fine-tuning
4 LLMs from OpenAI, Google, and DeepSeek with a small number of books unlocks the ability of those
5 models to reproduce up to 90% of text from unrelated copyrighted books, “using only semantic
6 descriptions as prompts and no actual book text.”⁷ The models identified in these studies rely on the
7 same or similar architecture as the Nemotron and SlimLM models developed by Adobe. Upon
8 information and belief, models developed by Adobe are capable of memorizing and therefore
9 regurgitating substantial portions of their training data.

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

16 27. Adobe, founded in 1982, provides services and tools for image, audio, video, and text
17 editing and viewing. It is well known as the creator of Photoshop (an image manipulation program),
18 and Adobe Acrobat (a program to view and manipulate files stored in Adobe’s proprietary Portable
19 Document Format (“PDF”) format).

20 28. Adobe has in recent years added artificial intelligence capabilities into its products.
21 Photoshop, for example, now incorporates a family of generative AI models known as Adobe Firefly
22 that are capable of generating and manipulating images, videos, and audio. Adobe claims that the
23 output of its Adobe Firefly models are “commercially safe” because Adobe used only licensed or public
24 domain images and videos to develop these models. According to Adobe, these images are “vetted and
25

26 ⁶ Ahmed Ahmed et. al., *Extracting Books from Production Language Models* (Jan. 6, 2026)
27 [<https://arxiv.org/pdf/2601.02671>].

28 ⁷ Xinyue Liu et al., *Alignment Whack-a-Mole: Finetuning Activates Verbatim Recall of Copyrighted
Books in Large Language Models* (2026) [<https://arxiv.org/abs/2603.20957>].

1 humanly verified to be copyright compliant,”⁸ and “customers can use those commercially safe
2 solutions to confidently publish outputs knowing Adobe has *responsibly addressed licensing and*
3 *copyright issues.*”⁹

4 29. Adobe Chief Strategy Officer Scott Belsky described other models in the industry as
5 built on data that is “openly scraped,”¹⁰ unlike Adobe’s AI models that “show respect for the creative
6 community” by training only on licensed content or public domain content where the copyright has
7 expired.¹¹ Adobe CEO Shantanu Narayen further noted:

8 I think other companies . . . are not completely transparent yet about what
9 data they use and [if] they scrape the internet, and that will play out in the
10 industry. But I like the approach that we’ve taken, and I like the way in
11 which we’ve engaged with our community on this.¹²

12 30. Adobe claims to purchase or acquire images for which “Adobe has monetarily
13 compensated the creator, secured the consent of person photographed, or confirmed that the copyright
14 license has expired.”¹³

15 31. Adobe even pays the creators of licensed works a yearly “Firefly Contributor Bonus”
16 based on the number of images used to train the Firefly models.¹⁴

17 ⁸ *Adobe Generative AI Built for Business Solution Brief*, ADOBE (2024), [https://www.adobe.com/cc-](https://www.adobe.com/cc-shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-brief.pdf)
18 [shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-](https://www.adobe.com/cc-shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-brief.pdf)
19 [brief.pdf](https://www.adobe.com/cc-shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-brief.pdf) (on file with counsel).

20 ⁹ *Id.* (emphasis added).

21 ¹⁰ *Adobe’s ‘Ethical’ Firefly AI Was Trained on Midjourney Images*, CNBC TV18 (Apr. 12, 2024),
22 [https://www.cnbc.com/technology/adobes-ethical-firefly-ai-was-trained-on-midjourney-images-](https://www.cnbc.com/technology/adobes-ethical-firefly-ai-was-trained-on-midjourney-images-19396150.htm)
23 [19396150.htm](https://www.cnbc.com/technology/adobes-ethical-firefly-ai-was-trained-on-midjourney-images-19396150.htm) (on file with counsel).

24 ¹¹ *Adobe Firefly vs. DALL-E 3*, ADOBE, [https://www.adobe.com/products/firefly/discover/firefly-vs-](https://www.adobe.com/products/firefly/discover/firefly-vs-dalle.html)
25 [dalle.html](https://www.adobe.com/products/firefly/discover/firefly-vs-dalle.html) (on file with counsel) (last visited May 17, 2026).

26 ¹² Nilay Patel, *Why Adobe CEO Shantanu Narayen is confident we’ll all adapt to AI*, THE VERGE (May
27 13, 2024, 7:00 AM PDT), [https://www.theverge.com/24153956/adobe-shantanu-narayen-ai-firefly-](https://www.theverge.com/24153956/adobe-shantanu-narayen-ai-firefly-premiere-photoshop-pdf-creativity-commerce)
28 [premiere-photoshop-pdf-creativity-commerce](https://www.theverge.com/24153956/adobe-shantanu-narayen-ai-firefly-premiere-photoshop-pdf-creativity-commerce) (on file with counsel).

29 ¹³ *Adobe Generative AI Built for Business Solution Brief*, ADOBE (2024), [https://www.adobe.com/cc-](https://www.adobe.com/cc-shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-brief.pdf)
30 [shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-](https://www.adobe.com/cc-shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-brief.pdf)
31 [brief.pdf](https://www.adobe.com/cc-shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-brief.pdf) (on file with counsel).

32 ¹⁴ Matthew Smith, *Adobe Firefly + Adobe Stock*, ADOBE BLOG (Sept. 13, 2023)
33 <https://blog.adobe.com/en/publish/2023/09/13/adobe-firefly-adobe-stock> (on file with counsel).

1 32. Adobe has gone as far as creating infrastructure to allow image creators to exclude their
2 works from the training data used for Adobe’s Firefly models. This is done through a “do not train” tag
3 affixed to images so that a creator’s content will not be used to train AI models.¹⁵ According to Adobe
4 Chief Technology Officer Ely Greenfield, Adobe’s AI models already prevent the output of copyrighted
5 content, stating that if Adobe users ask an AI model to generate an image in the style of a particular
6 artist, “it won’t generate an image that is aping that person’s style . . . If someone wants to use your
7 style, you can actually sell a customer the right to use your style.”¹⁶

8 33. Adobe acknowledges the legal risk associated with training an AI model with
9 copyrighted content. Adobe’s white paper, “Generative AI Built for Business, Adobe’s Commitment to
10 Building AI Responsibly,” notes that “if an individual or a company uses an image that has been
11 generated by an AI model trained on copyrighted or otherwise legally protected content, that individual
12 or company can be subject to reputational damage and expensive lawsuits for the use of that protected
13 content.”¹⁷ Adobe General Counsel Dana Rao further noted:

14 When we looked at Firefly or generative models and said, “How are we
15 going to train it? We could train it off the web. We could scrape the web and
16 build that model, or we could try to be more thoughtful about how we train
17 that model given the potential copyright issues and the fact that we have
18 creative customers who are concerned about people training on their content
19 without their permission.” . . . Now, the good news about the choice we
20 made on copyright is that it respects our creative customers who are
21 concerned about this. And enterprise customers were very excited to know
22 that they’re going to be able to use an image generative model that doesn’t
23 have IP issues, doesn’t have brand issues, isn’t trained on unsafe content,
24 because all of that has been either not in the database at all in the beginning
25 because of the way we trained it or we can use content moderation to
26 address it before it even gets into Firefly.¹⁸

22 ¹⁵ Dawn Chmielewski & Stephen Nellis, *Adobe, Nvidia AI imagery systems aim to resolve copyright*
23 *questions*, REUTERS (Mar. 21, 2023, at 12:52 PM PDT), [https://www.reuters.com/technology/adobe-](https://www.reuters.com/technology/adobe-nvidia-ai-imagery-systems-aim-resolve-copyright-questions-2023-03-21/)
24 [nvidia-ai-imagery-systems-aim-resolve-copyright-questions-2023-03-21/](https://www.reuters.com/technology/adobe-nvidia-ai-imagery-systems-aim-resolve-copyright-questions-2023-03-21/) (on file with counsel).

24 ¹⁶ *Id.*

25 ¹⁷ *Adobe Generative AI Built for Business Solution Brief*, ADOBE (2024), [https://www.adobe.com/cc-](https://www.adobe.com/cc-shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-brief.pdf)
26 [shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-](https://www.adobe.com/cc-shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-brief.pdf)
27 [brief.pdf](https://www.adobe.com/cc-shared/assets/pdf/trust-center/ungated/whitepapers/corporate/adobe-gen-ai-built-for-business-solution-brief.pdf) (on file with counsel).

27 ¹⁸ Nilay Patel, *How Adobe is managing the AI copyright dilemma, with general counsel Dana Rao*, THE
28 VERGE, (Jan 9, 2024, 7:00 AM PST), [https://www.theverge.com/24027198/adobe-dana-rao-ai-copyright-](https://www.theverge.com/24027198/adobe-dana-rao-ai-copyright-fair-use-figma-acquisition-deal-decoder-interview)
[fair-use-figma-acquisition-deal-decoder-interview](https://www.theverge.com/24027198/adobe-dana-rao-ai-copyright-fair-use-figma-acquisition-deal-decoder-interview) (on file with counsel).

1 34. Yet while making these representations, Adobe was simultaneously training its LLMs on
2 *hundreds of thousands* of pirated books without permission, without payment, and without any attempt
3 to obtain licenses from rightsholders. This is far from “showing respect for the creative community.” In
4 fact, Adobe did not even attempt to seek *any* licenses for the use of copyrighted books. It simply took
5 the works.

6 35. In March 2023, Adobe announced a partnership with NVIDIA to “create the next
7 generation of generative AI models” with a focus on transparency.¹⁹

8 36. A year later, in March 2024, Adobe again announced it was working together with
9 NVIDIA to “train new NVIDIA LLMs.”²⁰ Specifically, using Adobe’s PDF Extract tool, Adobe was
10 “building datasets to train and tune next-generation NVIDIA AI foundation models, including NVIDIA
11 Nemotron LLMs.”²¹

12 37. NVIDIA’s Nemotron models are a family of AI models trained on books. When NVIDIA
13 announced the release of the Nemotron-4 15B model in February 2024, it did not identify the specific
14 datasets used to train the model.²² NVIDIA instead indicated generally that the model was trained with
15 “books.”²³ The entire training corpus consisted of 8 trillion tokens. 5.6 trillion (or seventy percent) of
16 those tokens were made up of “English natural language” data, of which 240.8 billion (or 4.6%) were
17 English-language books.²⁴ To put these figures in perspective, one token is about four characters, or 3/4

18
19 ¹⁹ *NVIDIA Brings Generative AI to World’s Enterprises With Cloud Services for Creating Large*
20 *Language and Visual Models*, NVIDIA (Mar. 21, 2023), [https://investor.nvidia.com/news/press-release-](https://investor.nvidia.com/news/press-release-details/2023/NVIDIA-Brings-Generative-AI-to-Worlds-Enterprises-With-Cloud-Services-for-Creating-Large-Language-and-Visual-Models/default.aspx)
21 [details/2023/NVIDIA-Brings-Generative-AI-to-Worlds-Enterprises-With-Cloud-Services-for-Creating-](https://investor.nvidia.com/news/press-release-details/2023/NVIDIA-Brings-Generative-AI-to-Worlds-Enterprises-With-Cloud-Services-for-Creating-Large-Language-and-Visual-Models/default.aspx)
22 [Large-Language-and-Visual-Models/default.aspx](https://investor.nvidia.com/news/press-release-details/2023/NVIDIA-Brings-Generative-AI-to-Worlds-Enterprises-With-Cloud-Services-for-Creating-Large-Language-and-Visual-Models/default.aspx) (on file with counsel).

23 ²⁰ Abhigyan Modi, *Adobe Partners with NVIDIA to Harness the Power of PDF Intelligence with Next-*
24 *Gen LLMs*, ADOBE BLOG (Mar. 18, 2024), [https://blog.adobe.com/en/publish/2024/03/18/adobe-](https://blog.adobe.com/en/publish/2024/03/18/adobe-partners-nvidia-harness-power-pdf-intelligence-next-gen-llms)
25 [partners-nvidia-harness-power-pdf-intelligence-next-gen-llms](https://blog.adobe.com/en/publish/2024/03/18/adobe-partners-nvidia-harness-power-pdf-intelligence-next-gen-llms) (on file with counsel).

26 ²¹ *Id.*

27 ²² Jupinder Parmar et al., *Nemotron-4 15B Technical Report*, NVIDIA, Working Paper (Feb. 27, 2024)
28 [<https://arxiv.org/pdf/2402.16819>]. NVIDIA had not always concealed the datasets it used to train its
models. When releasing previous models, NVIDIA disclosed the specific datasets it used to train those
models, which revealed NVIDIA’s use of the Books3 dataset containing illicit copies of copyrighted
books.

²³ *Id.*

²⁴ *Id.*

1 of a word, and 100 tokens is about 75 words.²⁵ An average-length novel of 250 to 300 pages containing
2 250 to 300 words per page is about 128,000 tokens.²⁶ A dataset consisting of 240.8 billion tokens of
3 books data, with an average book length of 128,000 tokens, would contain **1,881,250** books.

4 38. In June 2024, NVIDIA released the Nemotron-4 340B model that used the same 8
5 trillion tokens from the Nemotron-4 15B model, but was fine-tuned on an additional 1 trillion tokens.²⁷
6 Adobe and NVIDIA have never announced any agreement with book authors or publishers to license
7 millions of books to train any AI model. In fact, a licensing deal of that magnitude does not exist *with*
8 *any AI company*. Upon information and belief, obtaining such a large number of English-language
9 books to train the Nemotron models could only be possible by pirating copyrighted books from sources
10 such as Anna's Archive that already make millions of copyrighted books available for download for
11 free.

12 39. Adobe and NVIDIA intentionally concealed the source of books data used to train the
13 Nemotron models to prevent exposure of their conduct and the enforcement of copyrights by
14 rightsholders.

15 40. Adobe itself admitted to training and fine-tuning Nemotron models, and building
16 datasets to train the Nemotron models using its Adobe PDF Extract tool²⁸ (an AI-powered tool allowing
17 the extraction of content from PDFs into a format suitable for LLM training).²⁹ This facilitates using the
18 text within PDF files to train or fine-tune LLMs.

19
20 ²⁵ *What are tokens and how to count them?*, OPENAI HELP CENTER,
21 <https://help.openai.com/en/articles/4936856-what-are-tokens-and-how-to-count-them> (on file with
counsel) (last visited May 17, 2026).

22 ²⁶ *Everything You Need to Know About Tokens, Data Volumes and Processing in ChatGPT*, UINTENT
23 (Nov. 26, 2024), [https://www.uintent.com/case-studies-und-blog/everything-you-need-to-know-about-
token-chatgpt](https://www.uintent.com/case-studies-und-blog/everything-you-need-to-know-about-token-chatgpt) (on file with counsel).

24 ²⁷ *Nemotron-4 340B Technical Report*, NVIDIA, Working Paper (Aug. 5, 2024)
25 [<https://arxiv.org/pdf/2406.11704>].

26 ²⁸ Abhigyan Modi, *Adobe Partners with NVIDIA to Harness the Power of PDF Intelligence with Next-
Gen LLMs*, ADOBE BLOG (Mar. 18, 2024), [https://blog.adobe.com/en/publish/2024/03/18/adobe-
partners-nvidia-harness-power-pdf-intelligence-next-gen-llms](https://blog.adobe.com/en/publish/2024/03/18/adobe-partners-nvidia-harness-power-pdf-intelligence-next-gen-llms) (on file with counsel).

27 ²⁹ *PDF Extract API Overview*, ADOBE DEVELOPER, [https://developer.adobe.com/document-
services/docs/overview/pdf-extract-api/](https://developer.adobe.com/document-services/docs/overview/pdf-extract-api/) (on file with counsel) (last visited May 17, 2026).
28

1 41. Pirated books are commonly stored in PDF format. Anna’s Archive, for example, makes
2 available for download PDF versions of books, including Plaintiff’s works.

3 42. NVIDIA acquired millions of unauthorized copies of copyrighted books from the
4 shadow library Anna’s Archive in order to develop its LLMs. Upon information and belief, Adobe used
5 these millions of pirated books to train and fine-tune the Nemotron models, including through the use
6 of its PDF Extract tool.

7 43. Shadow libraries are illicit online repositories of copyrighted content that openly defy
8 copyright laws. Typically, they allow anyone to search for and download pirated books, including in
9 PDF format. In spite of enforcement actions to curtail the ability of shadow libraries to distribute large
10 collections of books, several shadow libraries have flourished and gained notoriety. These include (1)
11 Library Genesis (“LibGen”), a website repeatedly enjoined by federal courts for copyright infringement
12 in default proceedings and which has been designated a “notorious” repository of pirated works by the
13 United States Trade Representative; its founders have faced federal criminal charges, (2) Z-Library (aka
14 B-ok), a shadow library that originally began as a for-profit LibGen mirror charging a fee for expedited
15 downloading, and which was seized by law enforcement; its founders have been arrested and indicted
16 with criminal copyright infringement (and have since fled the country), and (3) Sci-hub, another
17 shadow library operating in default of multiple federal court judgements for copyright infringement.

18 44. The most active and notorious shadow library is “Anna’s Archive,” which first began as
19 “Pirate Library Mirror” (named this way because it “mirrored” all of the books from Z-Library). The
20 shadow library rebranded to “Anna’s Archive” in 2022 and soon made available *all* of the content of
21 other shadow libraries, including LibGen, Z-Library, Sci-Hub, and Internet Archive, and additional
22 books sourced from pirated libraries. Anna’s Archive makes *millions* of pirated books available for
23 download to the public. Books can be downloaded individually, or entire datasets with hundreds of
24 thousands or millions of books can be downloaded wholesale through the peer-to-peer BitTorrent
25 Protocol (also known as “torrenting”). When someone downloads a file through BitTorrent, they must
26 also “seed” the file (that is, distribute it) to others who are also downloading the file. As a result, users
27 of the BitTorrent protocol are engaged in both piracy *and* the distribution of illegal content.

28

1 45. The use of pirated books from Anna’s Archive, LibGen, Z-Library, and other shadow
2 libraries has been commonplace within the AI development community.³⁰ Major AI companies,
3 including OpenAI, Meta, Anthropic, and NVIDIA pirated books from Anna’s Archive, Library Genesis,
4 Z-Library, Sci-Hub, and/or Pirate Library Mirror.

5 46. Shadow libraries provide an easily accessible source of high-quality copyrighted
6 material—for free. As Anna’s Archive explained, “[i]t is well understood that LLMs thrive on high-
7 quality data. We have the largest collection of books, papers, magazines, etc. in the world, which are
8 some of the highest quality text sources.”³¹ Shadow libraries themselves acknowledge that the
9 explosion in piracy and patronage by LLM companies has saved shadow libraries from extinction. As
10 the creators of Anna’s Archive wrote:

11 Not too long ago, “shadow-libraries” were dying. Sci-Hub, the massive
12 illegal archive of academic papers, had stopped taking in new works, due to
13 lawsuits. “Z-Library”, the largest illegal library of books, saw its alleged
14 creators arrested on criminal copyright charges ***Then came AI. Virtually all major companies building LLMs contacted us to train on our data. Most (but not all!) US-based companies reconsidered once they realized the illegal nature of our work . . . We have given high-speed access to about 30 companies. Most are Chinese, though we’ve also worked with companies from the US.***³²

16 47. Adobe’s disregard for the rights of copyright holders did not stop there. In August 2025,
17 Adobe released its SlimLM series of “Small Language Models” (“SLMs”), a type of LLM.³³ These
18 models were trained on the “SlimPajama” dataset that contains the Books3 dataset.
19
20
21

22 ³⁰ See, e.g., Alex Reisner, *The Unbelievable Scale of AI’s Pirated-Books Problem*, THE ATLANTIC (Mar.
23 20, 2025), <https://www.theatlantic.com/technology/archive/2025/03/libgen-meta-openai/682093/> (on file
24 with counsel); *Bartz v. Anthropic PBC*, 787 F. Supp. 3d 1007, 1015 (N.D. Cal. 2025) (noting Anthropic’s
25 use of LibGen and Pirate Library Mirror to download millions of copyrighted books).

26 ³¹ *LLM Data*, ANNA’S ARCHIVE, <https://annas-archive.pk/llm> (on file with counsel) (last visited May 17,
27 2026).

28 ³² *Copyright reform is necessary for national security*, ANNA’S ARCHIVE BLOG (Jan. 31, 2025)
<https://annas-archive.pk/blog/ai-copyright.html> (on file with counsel) (emphasis added).

³³ Thang M. Pham et al., *SlimLM: An Efficient Small Language Model for On-Device Document Assistance*, Working Paper (Nov. 25, 2024) [<https://arxiv.org/pdf/2411.09944>].

1 48. The SlimPajama dataset was created by AI development company Cerebras AI in June
2 2023 as a smaller version of the popular dataset known as RedPajama.³⁴

3 49. RedPajama (also known as RedPajama-V1) is a dataset created by Together AI in April
4 2023 for use in training AI models. The dataset aimed to “create a set of leading open-source models
5 and to rigorously understand the ingredients that yield good performance.”³⁵ In fact, the RedPajama
6 dataset specifically replicated the dataset mixture used by Meta to train the then-state-of-the-art
7 LLaMA series of models, noting that “for each data slice, [Together AI] conduct[ed] careful data pre-
8 processing and filtering. . . to roughly match the number of tokens as reported by Meta AI in the
9 LLaMA paper.”³⁶

10 50. One of the subsets of RedPajama is the “Books” dataset, described by Together AI as
11 “[a] corpus of open books.” Since RedPajama mirrored the dataset used to create the LLaMA models,
12 more detail about the RedPajama-Books subset is found in the paper “LLaMA: Open and Efficient
13 Foundation Language Models:”

14 We include two book corpora in our training dataset: the Gutenberg
15 Project, which contains books that are in the public domain, and the
16 Books3 section of The Pile (Gao et al., 2020), a publicly available dataset
17 for training large language models.³⁷

18
19
20
21

³⁴ Daria Soboleva, *SlimPajama: A 627B token, cleaned and deduplicated version of RedPajama*,
22 CEREBRAS BLOG (June 9, 2023), [https://www.cerebras.ai/blog/slimpajama-a-627b-token-cleaned-and-](https://www.cerebras.ai/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama)
23 [deduplicated-version-of-redpajama](https://www.cerebras.ai/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama) (on file with counsel).

24 ³⁵ *Releasing 3B and 7B RedPajama-INCITE family of models including base, instruction-tuned & chat*
25 *models*, TOGETHERAI BLOG (May 5, 2023), <https://www.together.ai/blog/redpajama-models-v1> (on file
26 with counsel).

27 ³⁶ *Id.*; *RedPajama, a project to create leading open-source models, starts by reproducing LLaMA*
28 *training dataset of over 1.2 trillion tokens*, TOGETHERAI BLOG (Apr. 17, 2023),
<https://www.together.ai/blog/redpajama> (on file with counsel).

³⁷ Hugo Touvron et al, *LLaMA: Open and Efficient Foundation Language Models*, Meta AI Working
Paper (Feb. 27, 2023) [<https://arxiv.org/pdf/2302.13971>] (on file with counsel).

1 51. The Books3 dataset was created by Shawn Presser in October 2020. Books3 is included
2 within the “SlimPajama-Books” subset used to train Adobe models. The paper, “The Pile: An 800GB
3 Dataset of Diverse Text for Language Modeling,”³⁸ details the contents of the Books3 dataset:

4 Books3 is a dataset of books derived from a copy of the contents of the
5 Bibliotik private tracker . . . Bibliotik consists of a mix of fiction and
6 nonfiction books and is almost an order of magnitude larger than our next
7 largest book dataset (BookCorpus2). We included Bibliotik because books
8 are invaluable for long-range context modeling research and coherent
9 storytelling.³⁹

10 52. Bibliotik is a notorious shadow library similar to the other shadow libraries described
11 above, such as Library Genesis (aka LibGen), Z-Library (aka B-ok), Sci-Hub, and Anna’s Archive. Mr.
12 Presser has confirmed in public statements that Books3 represents “all of Bibliotik” and contains
13 approximately 196,640 books.

14 53. The RedPajama dataset is available for download from Hugging Face. Before October
15 2023, the Books subset containing Books3 was “downloaded from Huggingface [sic]” when a user ran
16 the script that automatically assembled the RedPajama dataset.⁴⁰ But in October 2023, the Books3
17 dataset was removed with a message that it “is defunct and no longer accessible due to reported
18 copyright infringement.”⁴¹ After the “Books3 dataset” was removed from Hugging Face, the
19 RedPajama dataset documentation also added a message that Books3 is defunct “due to reported
20
21
22

23 ³⁸ Leo Gao et al., *The Pile: An 800GB Dataset of Diverse Text for Language Modeling*, EleutherAI
24 Working Paper (Dec. 31, 2020) [<https://arxiv.org/pdf/2101.00027.pdf>].

25 ³⁹ *Id.* at 3–4.

26 ⁴⁰ *RedPajama-Data-1T*, Hugging Face Datasets,
27 <https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T>
28 [<https://web.archive.org/web/20230420075601/https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T>].

⁴¹ *the_pile_books3*, Hugging Face Datasets, https://huggingface.co/datasets/the_pile_books3
[https://web.archive.org/web/20231127101818/https://huggingface.co/datasets/the_pile_books3].

1 copyright infringement.”⁴² Together AI has also acknowledged it took down Books3 “due to copyright
2 issues.”⁴³

3 54. The SlimPajama dataset was previously available for download on HuggingFace,⁴⁴ but
4 was removed entirely sometime in 2026. The SlimPajama dataset downloaded and used by Adobe
5 contained the SlimPajama-Books dataset that included Books3.

6 55. Plaintiff’s copyrighted books listed in Exhibit A are among the works available at Anna’s
7 Archive and included in the Books3 dataset. Below, these books are referred to as the Asserted Works.

8 56. Adobe benefitted commercially from its use of pirated books by incorporating Nemotron
9 and SlimLM models into its programs and services, thereby attracting more customers and revenue. For
10 example, in March 2026 Adobe and NVIDIA announced a collaboration in which “Adobe will bring
11 NVIDIA Nemotron capabilities to Adobe Acrobat to further elevate the quality of AI output and
12 increase productivity for business professionals, consumers and enterprises.”⁴⁵ Adobe Acrobat is the
13 widely used program to view and manipulate PDFs. The announcement also noted that Nemotron
14 models would power Adobe’s “agentic workflows.”⁴⁶ And in April 2026 Adobe announced it was
15 embedding “NVIDIA Nemotron open models into the CX Enterprise Coworker platform.”⁴⁷

16
17 ⁴² *RedPajama-Data-1T*, Hugging Face Datasets, <https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T>
18 [<https://web.archive.org/web/20240510231649/https://huggingface.co/datasets/togethercomputer/RedPajama-Data-1T>].

19 ⁴³ Maurice Weber et al., *RedPajama: an Open Dataset for Training Large Language Models*, Working
20 Paper (Nov. 19, 2024) [<https://arxiv.org/pdf/2411.12372>].

21 ⁴⁴ *SlimPajama-627B*, Hugging Face Datasets, <https://huggingface.co/datasets/cerebras/SlimPajama-627B>
22 [<https://web.archive.org/web/20230609191207/https://huggingface.co/datasets/cerebras/SlimPajama-627B>].

23 ⁴⁵ *Adobe and NVIDIA Announce Strategic Partnership to Deliver the Next Generation of Firefly Models and Creative, Marketing and Agentic Workflows*, ADOBE NEWSROOM (Mar. 16, 2026, at 1:30 PM), <https://news.adobe.com/news/2026/03/adobe-and-nvidia-announce-strategic-partnership> (on file with counsel).

24 ⁴⁶ *Id.*

25 ⁴⁷ Dom Nicastro, *Nvidia CEO Jensen Huang Tells the SaaS World: Agentic Is Here. Adobe Is Listening*, CMSWIRE (Apr. 20, 2026), <https://www.cmswire.com/digital-experience/nvidia-ceo-jensen-huang-told-the-saas-world-agentic-is-here-adobe-was-listening/> (on file with counsel).

1 57. Adobe infringed on Plaintiff and Class members' copyrighted works on a massive scale.
2 Adobe acquired books originating from the Books3 subset of SlimPajama without authorization from
3 or compensation to their authors. Upon information and belief, Adobe also acquired books originating
4 from Anna's Archive without authorization from or compensation to their authors. Adobe then
5 continued copying and storing the datasets, and used them to develop and train the Nemotron and
6 SlimLM series of LLMs.

7 58. Upon information and belief, Adobe used Plaintiff and Class members' copyrighted
8 works to train other models, including non-public models, and to develop other tools, programs, or
9 services.

10 59. In a speech at Adobe Summit 2023, Adobe CEO Shantanu Narayen stated that Adobe
11 has "been very transparent about the training data" Adobe used, while "a lot of the other companies are
12 actually using data that potentially—they're scraping the internet, they're accessing data that they may
13 or may not have rights and license to."⁴⁸ Rather than focusing on the conduct of other companies,
14 however, Mr. Narayen accurately described Adobe's own conduct: Adobe accessed datasets of hundreds
15 of thousands of books for which it held no rights or license to copy.

16 60. Mr. Narayen also championed the right of creators to exclude their works from being
17 used to train AI models, stating, "We think that that's important because people are going to say, and
18 our creative community is going to say, 'I want do not track. So I don't want to be part of that training
19 data.' So we created that entire infrastructure to allow people to do that."⁴⁹ Adobe, however, never
20 implemented any infrastructure to allow *book authors* to exclude their works from the training data of
21 Adobe's AI models.

22
23
24
25
26
27 ⁴⁸ *ADBE.OQ – Adobe. Inc. Presents at Adobe Summit 2023*, Edited Transcript (Mar. 21, 2023),
<https://www.adobe.com/cc-shared/assets/investor-relations/pdfs/adbe-summit2023-qna.pdf> (on file with
28 counsel).

⁴⁹ *Id.*

1 otherwise used to develop any large language model, program, or tool
2 during the Class Period.

3 **Books3 subclass.** All legal or beneficial owners of copyrighted works
4 registered with the United States Copyright Office originating from
5 Books3 that Adobe acquired, copied, stored, or otherwise used to
6 develop any large language model, program, or tool during the Class
7 Period

8 64. This Class definition excludes:

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

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

21 66. **Typicality.** Plaintiff's claims are typical of the claims of other members of the Class
22 because Plaintiff and all members of the Class were damaged by the same wrongful conduct of
23 Defendant as alleged herein, and the relief sought herein is common to all members of the Class.

24 67. **Adequacy.** Plaintiff will fairly and adequately represent the interests of the members of
25 the Class because the Plaintiff has experienced the same harms as the members of the Class and has no
26 conflicts with any other members of the Class. Furthermore, Plaintiff has retained sophisticated and
27 competent counsel who are experienced in prosecuting federal and state class actions, as well as other
28 complex litigation.

1 74. To develop its SlimLM language models, Adobe copied datasets containing Books3 and
2 incorporated them into the training dataset for its LLMs. Adobe made multiple copies of these datasets
3 during the development of the SlimLM series of LLMs.

4 75. Upon information and belief, to develop the Nemotron language models, Adobe copied
5 datasets originating from Anna’s Archive to develop these models and incorporate the datasets into the
6 models’ training data. Adobe made multiple copies of these datasets during the development of the
7 Nemotron models.

8 76. Plaintiff and the Class members never authorized Defendant to make copies of their
9 Asserted Works, make derivative works, publicly display copies (or derivative works), or distribute
10 copies (or derivative works). All those rights belong exclusively to Plaintiff and the Class members
11 under the U.S. Copyright Act.

12 77. By copying, storing, processing, reproducing, and using the datasets containing copies
13 of Plaintiff’s original works, Defendant has directly infringed Plaintiff’s exclusive rights in her
14 copyrighted works, in violation of the Copyright Act.

15 78. By copying, storing, processing, and reproducing the Nemotron and SlimLM models—
16 each trained on Plaintiff’s copyrighted works—Adobe has further directly infringed upon Plaintiff’s
17 exclusive rights in her copyrighted works, in violation of the Copyright Act.

18 79. Defendant repeatedly copied, stored, and used the Asserted Works without her
19 authorization or consent. Such copies were made in direct violation of her exclusive rights under the
20 Copyright Act.

21 80. Defendant’s infringing conduct, as alleged herein, was and continues to be willful and
22 carried out with full knowledge of Plaintiff’s rights in and to the Asserted Works. As a direct and
23 proximate result of its conduct, Defendant has wrongfully profited from copyrighted works to which it
24 holds no ownership interest.

25 81. By and through the actions alleged above, Defendant has infringed and will continue to
26 infringe Plaintiff’s copyrights.

1 82. Plaintiff has been injured by Defendant’s acts of direct copyright infringement. Plaintiff
2 is entitled to statutory damages, actual damages, restitution of profits, and all appropriate legal and
3 equitable relief.

4 **DEMAND FOR JUDGMENT**

5 Wherefore, Plaintiff requests that the Court enter judgment on their behalf and on behalf of the
6 Class defined herein, by ordering:

- 7 a) This action may proceed as a class action, with Plaintiff serving as Class Representative, and
8 with Plaintiff’s counsel as Class Counsel.
- 9 b) Judgment in favor of Plaintiff and the Class and against Defendant.
- 10 c) An award of statutory and other damages under 17 U.S.C. § 504 for Defendant’s violations
11 of the copyrights of Plaintiff and the Class.
- 12 d) An award of reasonable attorneys’ fees and reimbursement of all costs of this action
13 pursuant to 17 U.S.C. § 505, or as otherwise permitted by law.
- 14 e) A declaration that such infringement is willful.
- 15 f) An order directing the destruction or other reasonable disposition of all copies Defendant
16 made or used in violation of the exclusive rights of Plaintiff and the Class, pursuant to 17
17 U.S.C. § 503(b).
- 18 g) Injunctive relief sufficient to alleviate and stop Defendant’s unlawful conduct alleged herein.
- 19 h) Pre- and post-judgment interest on the damages awarded to Plaintiff and the Class, and that
20 such interest be awarded at the highest legal rate from and after the date this complaint is
21 first served on Defendant.
- 22 i) Defendant is responsible financially for the costs and expenses of a Court-approved notice
23 program through post and media designed to give immediate notification to the Class.
- 24 j) Further relief for Plaintiff and the Class as may be appropriate.

25 //

26 //

27 //

JURY TRIAL DEMANDED

Under Federal Rule of Civil Procedure 38(b), Plaintiff demands a trial by jury of all the claims asserted in this complaint so triable.

Dated: May 18, 2026

By: /s/ Joseph R. Saveri

Joseph R. Saveri (SBN 130064)
Diane S. Rice (SBN 118303) Christopher
K.L. Young (SBN 318371) William W.
Castillo Guardado (SBN 294159) **SAVERI
LAW FIRM, LLP**
550 California Street, Suite 910
San Francisco, CA 94104
Telephone: (415) 500-6800
Facsimile: (415) 395-9940
Email: jsaveri@saverilawfirm.com
drice@saverilawfirm.com
cyoung@saverilawfirm.com
wcastillo@saverilawfirm.com

*Attorneys for Individual and Representative Plaintiff
E. Molly Tanzer*