

1 [Counsel on signature page]  
2  
3  
4  
5  
6  
7

8 **UNITED STATES DISTRICT COURT**  
9 **NORTHERN DISTRICT OF CALIFORNIA**  
10 **OAKLAND DIVISION**

11 ABDI NAZEMIAN, et al,  
12 Individual and Representative Plaintiffs,  
13 vs.  
14 NVIDIA CORPORATION  
15 Defendant.

Master File Case No. 4:24-cv-01454-JST (SK)  
Consolidated with Case No. 4:24-cv-02655-JST  
(SK)

**FIRST CONSOLIDATED AMENDED  
COMPLAINT**

**Class Action**

**Demand for Jury Trial**

1 Plaintiffs Abdi Nazemian, Brian Keene, Stewart O’Nan, Andre Dubus III, and Susan Orlean  
 2 (together “Plaintiffs”), on behalf of themselves and all others similarly situated, bring this class  
 3 action complaint (“Complaint”) against Defendant NVIDIA Corporation (“NVIDIA” or  
 4 “Defendant”).

5 **OVERVIEW**

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

8 2. A *large language model* is an AI software program designed to emit convincingly  
 9 naturalistic text outputs in response to user prompts.

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

18 4. Plaintiffs and Class members are authors. They own registered copyrights in certain  
 19 books that NVIDIA has admitted copying, storing, and using to develop its AI language models.

20 5. NVIDIA copied these copyrighted works multiple times to train its language  
 21 models, including from known pirated libraries (also known as “shadow libraries”). Those notorious  
 22 shadow libraries include The Pile, Bibliotik, and Anna’s Archive.

23 I am on the data strategy team at NVIDIA, we are exploring including Anna’s  
 24 Archive in pre-training data for our LLMs.

25 We are figuring out internally whether we are willing to accept the risk of using  
 26 this data, but would like to speak with your team to get a better understanding of  
 27 LLM-related work you have done.

6. NVIDIA “got the green light” to use Anna’s Archive. NVIDIA did not hesitate in using pirated books from these illicit sources of copyrighted material, regardless of the “risk” or the harm to authors like the Plaintiffs.

7. And NVIDIA also caused numerous third parties to download and store Plaintiffs' copyrighted works by encouraging, facilitating, and promoting its customers to download copies of The Pile dataset, which includes more than one hundred thousand copyrighted books.

## **JURISDICTION AND VENUE**

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

9. Jurisdiction and venue are proper in this judicial district under 28 U.S.C. § 1331 because NVIDIA is headquartered in this district. NVIDIA created various large language models, including the NeMo Megatron models, and distributes them commercially. Therefore, a substantial part of the events giving rise to the claim occurred in this District. A substantial portion of the affected interstate trade and commerce was carried out in this District. Defendant has transacted business, maintained substantial contacts, and/or committed overt acts in furtherance of the illegal scheme and conspiracy throughout the United States, including in this District. Defendant's conduct has had the intended and foreseeable effect of causing injury to persons residing in, located in, or doing business throughout the United States, including in this District.

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

## PLAINTIFFS

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

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

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

14. Plaintiff Andre Dubus III is an author who lives in Massachusetts. Plaintiff Dubus owns registered copyrights in multiple books, including, *The Garden of Last Days*, *The Cage Keeper*, and *Townie: A Memoir*.

15. Plaintiff Susan Orlean is an author who lives in California. Plaintiff Orlean owns registered copyrights in multiple works, including, *The Orchid Thief* and *The Library Book*.

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

**DEFENDANT**

17. Defendant NVIDIA is a Delaware corporation with its principal place of business at 2788 San Tomas Expressway, Santa Clara CA 95051.

## AGENTS AND CO-CONSPIRATORS

18. The unlawful acts alleged against the Defendant in this class action 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 businesses 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.

19. 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

20. NVIDIA is a diversified technology company founded in 1993 that originally focused on computer-graphics hardware, e.g., Graphics Processing Units (“GPUs”), and has since expanded to other computationally intensive fields, including software such as NVIDIA’s “Compute Unified Device Architecture” and hardware, e.g. NVLink/NVLink Switch, for training and operating AI software programs. NVIDIA’s hardware and software is used by all Frontier AI companies—companies that develop the most advanced AI systems—which has resulted in NVIDIA becoming the world’s most valuable company.

21. In addition to the hardware and software products it sells to AI companies, NVIDIA itself has developed numerous AI models known as “large language models” (“LLMs”). An LLM is AI software designed to emit convincingly naturalistic text outputs in response to user prompts. NVIDIA sells products to its customers that rely on NVIDIA’s LLMs.

22. Though LLMs are software programs, they are not created the way most software programs are—that is, by human software programmers writing code. Rather, LLMs are *trained* by copying an enormous quantity of textual works and then feeding these copies in pieces into the model. This corpus of input material is called the *training dataset*.

23. As set forth below, NVIDIA unlawfully copied copyrighted material from illegal pirate “shadow libraries.” NVIDIA collated and stored this material in centralized servers which its engineers (and other employees) could access for any purpose. NVIDIA and its employees subsequently made additional unlawful copies of this illegally-obtained copyrighted material during the LLM development process.

24. During the training process, LLMs copy and ingest each textual work in the training dataset and extract protected expression from it. In a process somewhat resembling a guess-and-check quiz, the LLM is progressively adjusted to more closely approximate the protected expression copied from the training dataset. The LLM records the results of this process in a large set of numbers called *weights* or *parameters* that are stored within the model, and, in some sense, “are” the model. These weights are entirely and uniquely derived from the protected expression in the training dataset.

1 For instance, the NeMo Megatron–GPT 20B model—an LLM released in September 2022 as part  
 2 of NVIDIA’s NeMo Megatron series of LLMs—is so named because the model stores 20 billion  
 3 (“20B”) weights derived from protected expression in its training dataset.

4 25. Importantly, datasets may have multiple uses during the development process of an  
 5 LLM even if the dataset does not become part of a model’s final training dataset. For example, during  
 6 the development of an LLM, the developer may initiate a *run* or *checkpoint* using certain datasets to  
 7 see the effect of that dataset on the model. Once the checkpoint is finished, a full model is completed  
 8 and its performance analyzed. The developer may then alter the datasets and conduct another  
 9 checkpoint. This process may occur multiple times before a developer arrives at the final checkpoint  
 10 for that model. All of the models created as part of the checkpoint process may never receive official  
 11 names nor be publicly released.

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

17 27. Much of the material in NVIDIA’s training dataset, however, comes from  
 18 copyrighted works—including books written by Plaintiffs and Class members—that were acquired,  
 19 copied and stored by NVIDIA without consent, without credit, and without compensation.

20 28. In November 2021, NVIDIA announced the “NeMo Megatron framework for  
 21 training language models.”<sup>1</sup> NVIDIA touted this framework as “provid[ing] a production-ready,  
 22 enterprise-grade solution to simplify the development and deployment of large language models.”<sup>2</sup>

23 29. In September 2022, NVIDIA announced the availability of the NeMo Megatron  
 24 language models in a video on its website: “For the first time, NVIDIA is making its checkpoints  
 25 available publicly, where the checkpoints are trained with NeMo Megatron … this is just to begin

27 <sup>1</sup> See <https://nvidianews.nvidia.com/news/nvidia-brings-large-language-ai-models-to-enterprises-worldwide>.

28 <sup>2</sup> *Id.*

1 with. And this is not the end. We will continue to add more checkpoints in the future.”<sup>3</sup> In this context  
 2 “checkpoints” is an alternate term for language models. The language models released in September  
 3 2022 include NeMo Megatron-GPT 1.3B, NeMo Megatron-GPT 5B, NeMo Megatron-GPT 20B,  
 4 and NeMo Megatron-T5 3B models.

5       30.     Each of these NeMo Megatron models was hosted on a website called Hugging  
 6 Face, where a model card provides information about each model, including its training dataset. The  
 7 model card for each of the NeMo Megatron models states, “The model was trained on ‘The Pile’  
 8 dataset prepared by EleutherAI.”<sup>4</sup>

9       31.     The Pile is a training dataset curated by a research organization called EleutherAI.  
 10 In December 2020, EleutherAI introduced this dataset in a paper called “The Pile: An 800GB Dataset  
 11 of Diverse Text for Language Modeling”<sup>5</sup> (the “EleutherAI Paper”).

12       32.     According to the EleutherAI Paper, one of the components of The Pile is a collection  
 13 of books called Books3. The EleutherAI Paper reveals that the Books3 dataset comprises  
 14 108 gigabytes of data, or approximately 12% of the dataset, making it the third largest component  
 15 of The Pile by size.

16       33.     The EleutherAI Paper further describes the contents of Books3:

17                 Books3 is a dataset of books derived from a copy of the contents of  
 18 the Bibliotik private tracker … Bibliotik consists of a mix of fiction  
 19 and nonfiction books and is almost an order of magnitude larger than  
 20 our next largest book dataset (BookCorpus2). **We included Bibliotik**  
 21 **because books are invaluable for long-range context modeling**  
 22 **research and coherent storytelling.**<sup>6</sup>

23

---

24       <sup>3</sup> See <https://www.nvidia.com/en-us/on-demand/session/gtcfall22-a41200/?nvid=nv-int-tblg-881125>, starting at 37:25.

25       <sup>4</sup> See, e.g., <https://huggingface.co/nvidia/nemo-megatron-gpt-1.3B#training-data>,  
 26 <https://huggingface.co/nvidia/nemo-megatron-gpt-5B#training-data>,  
<https://huggingface.co/nvidia/nemo-megatron-gpt-20B#training-data>,  
<https://huggingface.co/nvidia/nemo-megatron-t5-3B#training-data>

27       <sup>5</sup> Available at <https://arxiv.org/pdf/2101.00027.pdf>

28       <sup>6</sup> *Id.* at 3–4 (emphasis added).

1       34.     Bibliotik is one of a number of notorious “shadow library” websites which make,  
 2 store, and distribute huge quantities of pirated copyrighted works via the BitTorrent Protocol.

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

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

7       37.     Until October 2023, the Books3 dataset was available from Hugging Face. At that  
 8 time, the Books3 dataset was removed with a message that it “is defunct and no longer accessible  
 9 due to reported copyright infringement.”<sup>7</sup>

10      38.     NVIDIA has publicly admitted training its NeMo Megatron models on a copy of  
 11 The Pile dataset. Therefore, NVIDIA necessarily also (1) acquired a copy of Books3 (because it is  
 12 part of The Pile) and (2) made additional copies of Books3 during the course of developing LLMs,  
 13 including (but not limited to) its NeMo Megatron models. Certain books written by Plaintiffs are  
 14 part of Books3—including the Infringed Works—and thus NVIDIA necessarily (1) made unlawful  
 15 copies of Plaintiffs’ works when downloading Books3, and (2) made additional unlawful copies of  
 16 Plaintiffs’ works when developing its LLMs, including (but not limited to) its NeMo Megatron  
 17 models. NVIDIA thus directly infringed Plaintiffs’ copyrights.

18      39.     But NVIDIA’s use of Plaintiffs’ Infringed Works was not limited to the models it  
 19 *publicly* disclosed were trained on The Pile. NVIDIA and its engineers maintained The Pile in  
 20 centralized servers and repeatedly (and extensively) used The Pile following its acquisition,  
 21 including to develop multiple LLMs known internally as NeMo Megatron GPT 126M, NeMo  
 22 Megatron GPT 40B, NeMo Megatron GPT 175B, NeMo Megatron T5 220M, NeMo Megatron T5  
 23 11B, and NeMo Megatron T5 23B.

24      40.     NVIDIA’s use of The Pile to develop language models was not limited to a single  
 25 line or class of models either. Instead, language models across NVIDIA used The Pile.

26  
 27  
 28      <sup>7</sup> See [https://huggingface.co/datasets/the\\_pile\\_books3](https://huggingface.co/datasets/the_pile_books3)

1       41.     NVIDIA used The Pile to train and develop models that do not bear the NeMo  
 2 Megatron name as well. For instance, NVIDIA included the Pile dataset as training data for an LLM  
 3 known as Megatron 345M, which was publicly released as the Megatron GPT2 345m model.  
 4 NVIDIA also used The Pile to train an LLM known as “NeMo GPT-3 10B.” NVIDIA additionally  
 5 developed the InstructRetro-48B and Retro-48B LLMs using the Books3 dataset from The Pile.

6       42.     The Pile was not NVIDIA’s only dataset that included Books3. NVIDIA also  
 7 downloaded the SlimPajama dataset.<sup>8</sup> “SlimPajama was created by cleaning and deduplicating the  
 8 1.2T token RedPajama dataset from [the company] Together [AI].” And the RedPajama dataset itself  
 9 originally included the Books3 dataset. The SlimPajama dataset included the Books3 dataset.  
 10 NVIDIA used the SlimPajama dataset to test “both sentencepiece and BPE [tokenizers].” Tokenizers  
 11 are software which is used to process training data for use in LLM training and development. In  
 12 short, NVIDIA used the SlimPajama dataset to develop and test the software used in the development  
 13 of its LLMs. As one NVIDIA employee remarked, “SlimPajama . . . is available in our org.”  
 14 NVIDIA, therefore, again infringed Plaintiffs’ copyrights by downloading unauthorized copies of  
 15 their works by downloading, storing, and using the SlimPajama dataset.

16       43.     Upon information and belief, NVIDIA also developed a large number of internal  
 17 models, including checkpoints, many of which were never given proper names or publications but  
 18 which also unlawfully included datasets containing Plaintiffs’ and Class members’ works, such as  
 19 The Pile.

20       44.     Upon information and belief, NVIDIA also made unlawful copies of The Pile during  
 21 the course of internal research which did not result in a fully trained LLM.

22       45.     Not content to acquire, store, and use The Pile in its internal and external LLM  
 23 research, development, and commercialization efforts, NVIDIA sought vastly more copyrighted  
 24 works than The Pile could provide. Because the quality of an LLM depends on both the quality *and*  
 25 quantity of its training data, NVIDIA found itself desperate for additional books. Books have the  
 26  
 27

28       <sup>8</sup> See <https://huggingface.co/datasets/cerebras/SlimPajama-627B>.

1 unique designation of being widely understood as high-quality LLM training data and being  
 2 available illegally in large quantities from illicit shadow libraries.

3       46.      In addition to Bibliotik (the source of Books3, discussed above), those shadow  
 4 libraries include: (1) Library Genesis (“LibGen”) which has been repeatedly enjoined by federal  
 5 courts for copyright infringement in default proceedings and which has been designated a  
 6 “notorious” repository of pirated works by the United States Trade Representative; (2) Z-Library  
 7 (aka B-ok) which began as a for-profit LibGen mirror which enabled expedited downloads for a fee  
 8 until it was seized by law enforcement as part of an operation which resulted in its founders being  
 9 arrested and indicted (they have since fled the country); and (3) Sci-hub which, like LibGen, has  
 10 been repeatedly enjoined by federal courts for copyright infringement in default proceedings.

11       47.      The most active current shadow library is known as “Anna’s Archive.” The successor  
 12 to Z-library, Anna’s Archive began existence as “Pirate Library Mirror,” a name derived from the  
 13 fact that it “mirrored” (that is to say, hosted all the same books as) Z-Library. Shortly after its launch  
 14 in 2022, it rebranded to “Anna’s Archive” and quickly expanded to host all of LibGen, Z-Library,  
 15 Sci-Hub, and additional books sourced from pirated libraries. Anna’s Archive hosts millions of  
 16 pirated books.

17       48.      Many of these shadow libraries enable increased download speeds or quantities for  
 18 paying members. *See, e.g.*, <https://annas-archive.org/donate>.

19       49.      These “shadow libraries” have long been of interest to the AI industry—and their  
 20 insatiable quest for more data—because they illegally host and distribute vast quantities of high-  
 21 quality copyrighted material and because they are willing to move LLM developers to the “front of  
 22 the line” for download speeds—in exchange for a fee.

23       50.      As Anna’s Archive explained, “[i]t is well understood that LLMs thrive on high-  
 24 quality data. We have the largest collection of books, papers, magazines, etc. in the world, which are  
 25 some of the highest quality text sources.” <https://annas-archive.org/llm>. Shadow libraries provide  
 26 “high-speed . . . enterprise-level access [to their collections] . . . [in exchange] for donations in the  
 27 range of tens of thousands USD.” In other words: paid piracy.

1       51.       As revealed publicly over the last year,<sup>9</sup> it is an industry-wide practice to use shadow  
 2 libraries such as Library Genesis, Z-Library, and Pirate Library Mirror. Virtually every one of the  
 3 major LLM developers—including OpenAI, Meta, and Anthropic—pirated books from Library  
 4 Genesis, Z-Library, Sci-Hub, and/or Pirate Library Mirror. NVIDIA followed this industry-wide  
 5 practice and pirated troves of books from shadow libraries.

6       52.       The shadow libraries themselves have noted that the explosion in piracy and  
 7 patronage by LLM companies has saved shadow libraries from extinction. As a post by the admins  
 8 of Anna’s Archive put it:

9               Not too long ago, “shadow-libraries” were dying. Sci-Hub, the  
 10 massive illegal archive of academic papers, had stopped taking in  
 11 new works, due to lawsuits. “Z-Library”, the largest illegal library of  
 12 books, saw its alleged creators arrested on criminal copyright charges  
 13 . . . *Then came AI. Virtually all major companies building LLMs*  
 14 *contacted us to train on our data. . . We have given high-speed*  
 15 *access to about 30 companies.* <https://annas-archive.org/blog/ai-copyright.html> (emphasis added.)

16       53.       Internal documents show competitive pressures drove NVIDIA to piracy. In the fall  
 17 of 2023, NVIDIA faced a rapidly approaching deadline in the form of its annual developer day. In  
 18 the year since the launch of the NeMo Megatron series in September 2022, OpenAI had released  
 19 ChatGPT to massive success, resulting in a substantial increase in investor attention on AI. In  
 20 response, NVIDIA sought to develop and demonstrate cutting edge LLMs at its fall 2023 developer  
 21 day. In seeking to acquire data for what it internally called “NextLargeLLM,” “NextLLMLarge” and  
 22 “Next Generation LLM” (collectively, “NextLargeLLM”). NVIDIA was “[h]yper [f]ocused on  
 23 books corpuses.” NVIDIA knew that “published books under copyright” are “the most valuable” for  
 24 developing LLMs and NVIDIA knew that only books were available in sufficient quantities. And

---

26       9       See, e.g., Alex Reisner, *The Unbelievable Scale of AI’s Pirated-Books Problem*, The Atlantic  
 27 (March 20, 2025), <https://www.theatlantic.com/technology/archive/2025/03/libgen-meta-openai/682093/>; *Bartz v. Anthropic PBC*, 787 F. Supp. 3d 1007, 1015 (N.D. Cal. 2025) (noting  
 28 Anthropic’s use of LibGen and Pirate Library Mirror to download millions of copyrighted books).

1 NVIDIA needed to achieve 8 trillion tokens for the “NextLargeLLM,” and books provided this  
 2 means.

3       54.     In August 2023, NVIDIA contacted books publishers to obtain fast “access to large  
 4 volumes of unique, high-quality datasets” or “ie. books.” But on information and belief, NVIDIA  
 5 could not secure this fast access to the huge quantity of books it needed through publishers. As one  
 6 book publisher told NVIDIA, it was “ not in a position to engage directly just yet but will be in  
 7 touch.” In 2023, NVIDIA had “chatted with multiple publishers . . . but none [] wanted to enter into  
 8 data licensing deals.”

9       55.     Desperate for books, NVIDIA contacted Anna’s Archive—the largest and most  
 10 brazen of the remaining shadow libraries—about acquiring its millions of pirated materials and  
 11 “including Anna’s Archive in pre-training data for our LLMs.” Because Anna’s Archive charged tens  
 12 of thousands of dollars for “high-speed access” to its pirated collections, *see <https://annas-archive.org/llm>*, NVIDIA sought to find out what “high-speed access” to the data would look like.

14       56.     In correspondence with NVIDIA executives, Anna’s Archive stated that, because its  
 15 collections were illegally acquired and maintained, NVIDIA executives would need to “let [Anna’s  
 16 Archive] know when you have decided internally that this is something you can pursue. We have  
 17 wasted too much time on people who could not get internal buy-in.”

18       57.     Within a week of contacting Anna’s Archive, and days after being warned by Anna’s  
 19 Archive of the illegal nature of their collections, NVIDIA management gave “the green light” to  
 20 proceed with the piracy. Anna’s Archive offered NVIDIA millions of pirated copyrighted books.  
 21 Anna’s Archive also offered access to several million books from Internet Archive, which were only  
 22 normally available through Internet Archive’s digital lending system (a system which was found to  
 23 be copyright infringement by the Second Circuit, *see Hachette Book Grp., Inc. v. Internet Archive*,  
 24 115 F.4th 163 (2d Cir. 2024)). Anna’s Archive promised NVIDIA access to “a lot of books,” totaling  
 25 roughly 500 terabytes of data. By downloading Anna’s Archive, NVIDIA pirated additional copies  
 26 of Plaintiff’s Infringed Works.

1       58. On information and belief, in addition to Anna’s Archive and The Pile, NVIDIA also  
 2 downloaded books hosted or sourced from other shadow libraries, including LibGen, Sci-Hub, and  
 3 Z-Library.

4       59. About four months after its exchange with Anna’s Archive, in February 2024,  
 5 NVIDIA released a model known as Nemotron-4 15B. The training data for this model was not  
 6 publicly disclosed. Public documents, however, indicate that it was trained on 8 trillion tokens. The  
 7 sources of the training data were never identified, and NVIDIA stated that it included “books.”  
 8 NVIDIA, however, has publicly stated that the training data for this model encompasses 70% from  
 9 an “English natural language” dataset. This dataset itself is composed of 4.6% of books. Upon  
 10 information and belief, to reach this percentage of tokens derived from books, the training data  
 11 would need to include millions of books.

12       60. And a few months later, NVIDIA released the Nemotron-4 340B model. This model  
 13 included the same 8 trillion tokens from the Nemotron-4 15B but added an additional 1 trillion  
 14 tokens.

15       61. Upon information and belief, NVIDIA could not obtain the level of books needed  
 16 for the Nemotron models without pirating copyrighted books, including Plaintiffs’ Infringed Works.

17       62. In sum, NVIDIA has extensively and repeatedly violated the copyrights of Plaintiffs’  
 18 Infringed Works including by acquiring these works from pirated sources, storing them, and  
 19 enabling its employees to use them for any purpose, and copying them during the LLM training  
 20 process.

21       63. Plaintiff Abdi Nazemian’s book, *Like a Love Story*, was included in the Books3  
 22 dataset, based on public reporting about the dataset. This work is also available online through  
 23 Anna’s Archive, LibGen, and Z-Library.

24       64. Plaintiff Brian Keene’s book, *Ghost Walk*, was included in the Books3 dataset, based  
 25 on public reporting about the dataset. This work is also available online through Anna’s Archive,  
 26 LibGen, Z-Library, and Internet Archive.

1       65. Plaintiff Stewart O’Nan’s book, *Last Night at the Lobster*, was included in the  
 2 Books3 dataset, based on public reporting about the dataset. This work is also available online  
 3 through Anna’s Archive, LibGen, Z-Library, and Internet Archive.

4       66. Plaintiff Andre Dubus’s books, *The Garden of Last Days*, *The Cage Keeper*, and  
 5 *Townie: A Memoir* were included in the Books3 dataset, based on public reporting about the dataset.  
 6 These works are also available online through Anna’s Archive, LibGen, Z-Library, and Internet  
 7 Archive.

8       67. Plaintiff Susan Orlean’s books, *The Orchid Thief* and *The Library Book* were  
 9 included in the Books3 dataset, based on public reporting about the dataset. These works are also  
 10 available online through Anna’s Archive, LibGen, and Z-Library.

11       68. NVIDIA’s infringing activities, however, were not limited to downloading pirated  
 12 copyrighted material to develop and train its own language models. NVIDIA also provided the tools  
 13 and means for numerous others to infringe Plaintiffs’ copyrights.

14       69. As CEO Jensen Huang explained in the keynote address at NVIDIA’s 2023 GPU  
 15 Technology Conference, as part of NVIDIA’s “AI Foundations,” customers can use the NeMo  
 16 Framework (otherwise known as the NeMo Megatron Framework), to create and build their own AI  
 17 models. As he stated, “[t]hroughout the entire process, NVIDIA AI experts will work with you, from  
 18 creating your proprietary model to operations.”<sup>10</sup> As part of this process, NVIDIA assisted and  
 19 encouraged its customers to infringe Plaintiffs’ copyrights.

20       70. Through the NeMo Megatron Framework and BigNLP platforms, NVIDIA provided  
 21 customers with “scripts to automatically download and preprocess The Pile dataset which, until  
 22 recently, was hosted externally by Eleuther AI.” Meaning, NVIDIA provided tools and resources  
 23 for its customers to use the NVIDIA platform to download The Pile, thereby infringing on Plaintiffs’  
 24 copyrights. They scripts were developed to help their customers access these pirated datasets more  
 25 quickly and easily. NVIDIA employees expressed concern about the “[t]ime needed for downloading

26  
 27  
 28       <sup>10</sup> <https://www.youtube.com/watch?v=DiGB5uAYKAg> (40:00-:45).

pile files,” so they developed and distributed code to “download and extract[] 30 pile files [in] ~70 minutes[,] which clearly shows the need for data prep parallelism.”

71. For example, NVIDIA provided resources, guidance, and tools for its customer Writer Inc. to develop its line of Palmyra models using the NeMo Megatron Framework. On information and belief, NVIDIA provided the tools and scripts for Writer to download The Pile. NVIDIA provided similar assistance in downloading and processing The Pile to clients Persimmon AI Labs and Amazon. On information and belief, NVIDIA materially aided numerous other customers in downloading, using, and storing The Pile (and Books3) dataset.

72. NVIDIA provided the hardware too. Using the NeMo Framework, a customer could expect to quickly develop a language model trained on The Pile in only 9.8 days using NVIDIA's servers.

73. NVIDIA directly benefited from facilitating, supporting, and encouraging these infringing activities and attracted customers to use the NeMo Megatron Framework by providing quick access to The Pile (and Plaintiffs' books). In short, The Pile (and Books3) was key to NVIDIA attracting customers, and NVIDIA materially aided its customers to infringe Plaintiffs copyrights.

**COUNT 1**  
**Direct Copyright Infringement (17 U.S.C. § 501)**  
**against NVIDIA**

74. Plaintiffs incorporate by reference the preceding factual allegations.

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

76. To develop NVIDIA's LLMs, NVIDIA downloaded and copied The Pile and SlimPajama datasets. The Pile and SlimPajama datasets include the Books3 dataset, which includes the Infringed Works. NVIDIA made multiple copies of the Books3 dataset while developing its LLMs.

77. To develop NVIDIA's LLMs, NVIDIA downloaded and copied a dataset of books from Anna's Archive, which includes the Infringed Works. NVIDIA made multiple copies of this dataset while training its LLMs.

78. On information and belief, NVIDIA downloaded books hosted or sourced from other shadow libraries, including LibGen, Sci-Hub, and Z-Library.

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

80. NVIDIA made multiple copies of the Infringed Works, including when it downloaded these works from shadow libraries, and when it made additional copies during the training and development of its language models without Plaintiffs' permission and in violation of their exclusive rights under the Copyright Act. On information and belief, NVIDIA has continued to store and make copies of the Infringed Works.

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

82. NVIDIA's violation of Plaintiffs' and Class members' exclusive right was willful because NVIDIA knew the datasets it downloaded, copied, and stored, and on which it "trained" its LLMs contained copyrighted works.

**COUNT II**  
**Contributory Copyright Infringement**  
**against NVIDIA**

83. Plaintiffs incorporate by reference the preceding factual allegations.

84. NVIDIA materially contributed to and directly assisted in the direct infringement by multiple customers, including at least Amazon, Persimmon AI, and Writer, by providing the technology, personnel, access to datasets, and other resources, such as the NeMo Megatron Framework, and variations of similar platforms and scripts that performed the same function;

1 controlling or managing the property or other assets with which the direct infringement was  
 2 accomplished; or providing business, legal, strategic, or operational guidance that allowed its  
 3 customers to download, copy, and store Plaintiffs' and Class members' copyrighted works.

4 85. NVIDIA knew or had reason to know of the direct infringement by others using the  
 5 NeMo Megatron framework, because NVIDIA is fully aware of the capabilities of its own product,  
 6 platforms and tools upon which third parties downloaded and acquired at least The Pile dataset, and  
 7 potentially other datasets including copyrighted books as well.

8 86. Defendant is contributorily liable for the direct infringement of others that used the  
 9 NeMo Framework to download and acquire The Pile dataset (and potentially other datasets  
 10 containing copyrighted books as well).

11 **COUNT III**  
 12 **Vicarious Copyright Infringement**  
 13 **against NVIDIA**

14 87. Plaintiffs incorporate by reference the preceding factual allegations.

15 88. NVIDIA had the right and ability to control the direct infringements of customers,  
 16 including at least Amazon, Persimmon AI, and Writer, using the NeMo Megatron Framework, and  
 17 variations of similar platforms and scripts provided by NVIDIA that performed the same function,  
 18 to download The Pile dataset (and potentially other datasets containing copyrighted books as well).  
 NVIDIA failed to exert its right and ability to control its customers infringing acts.

19 89. NVIDIA has directly benefitted financially from the direct infringement of its  
 20 customers because NVIDIA generated revenue from customers using the NeMo Megatron  
 21 Framework to download The Pile (and Books3) dataset (and potentially other datasets containing  
 22 copyrighted books as well).

23 90. Plaintiffs have been injured by NVIDIA's acts of vicarious copyright infringement.  
 24 Plaintiffs are entitled to statutory damages, actual damages, restitution of profits, and other remedies  
 25 provided by law.

## CLASS ALLEGATIONS

91. The “**Class Period**” as defined in this Complaint begins no later than March 8, 2021 and runs through the present. Because Plaintiffs do not yet know when the unlawful conduct alleged herein began, but believe, on information and belief, that the conduct likely began earlier than March 8, 2021, Plaintiffs reserve the right to amend the Class Period to comport with the facts and evidence uncovered during further investigation or through discovery.

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

**All persons or entities that own a registered United States copyright in any literary work that was downloaded or otherwise copied by Defendant and / or used by Defendant in LLM training, research, or development during the Class Period.**

93. This Class definition excludes:

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

94. **Numerosity.** Plaintiffs do not know the exact number of members in the Class. This information is in the exclusive control of Defendant. On information and belief, there are at least tens or hundreds of thousands of members in the Class geographically dispersed throughout the United States. Therefore, joinder of all members of the Class in the prosecution of this action is impracticable.

1       95.     **Typicality.** Plaintiffs' claims are typical of the claims of other members of the Class  
 2 because Plaintiffs and all members of the Class were damaged by the same wrongful conduct of  
 3 Defendant as alleged herein, and the relief sought herein is common to all members of the Class.

4       96.     **Adequacy.** Plaintiffs will fairly and adequately represent the interests of the  
 5 members of the Class because the Plaintiffs have experienced the same harms as the members of the  
 6 Class and have no conflicts with any other members of the Class. Furthermore, Plaintiffs have  
 7 retained sophisticated and competent counsel who are experienced in prosecuting federal and state  
 8 class actions, as well as other complex litigation.

9       97.     **Commonality and predominance.** Numerous questions of law or fact common to  
 10 each Class member arise from Defendant's conduct and predominate over any questions affecting  
 11 the members of the Class individually:

- 12       a. Whether Defendant violated the copyrights of Plaintiffs and the Class when they  
 13            obtained copies of Plaintiffs' Infringed Works
- 14       b. Whether Defendant violated the copyrights of Plaintiffs and the Class when they used  
 15            them to research, develop, and train language models.
- 16       c. Whether Defendant intended to cause further infringement of the Infringed Works  
 17            with these language models because they have distributed these models under an  
 18            open license and advertised those models as a base from which to build further  
 19            models.
- 20       d. Whether Defendant's support, facilitation, and encouragement of the infringement  
 21            by NVIDIA's customers of Plaintiffs' and Proposed Class Members' copyrighted  
 22            works constitutes vicarious or contributory infringement under the Copyright Act
- 23       e. Whether any affirmative defense excuses Defendant's conduct.
- 24       f. Whether any statutes of limitation constrain the potential for recovery for Plaintiffs  
 25            and the Class.

26       98.     **Other class considerations.** Defendant has acted on grounds generally applicable  
 27 to the Class. This class action is superior to alternatives, if any, for the fair and efficient adjudication  
 28

1 of this controversy. Prosecuting the claims pleaded herein as a class action will eliminate the  
 2 possibility of repetitive litigation. There will be no material difficulty in the management of this  
 3 action as a class action. The prosecution of separate actions by individual Class members would  
 4 create the risk of inconsistent or varying adjudications, establishing incompatible standards of  
 5 conduct for Defendant.

6 **DEMAND FOR JUDGMENT**

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

- 9 a) This action may proceed as a class action, with Plaintiffs serving as Class  
     10 Representatives, and with Plaintiffs' counsel as Class Counsel.
- 11 b) Judgment in favor of Plaintiffs and the Class and against Defendant.
- 12 c) An award of statutory and other damages under 17 U.S.C. § 504 for violations of  
     13 the copyrights of Plaintiffs and the Class by Defendant.
- 14 d) Reasonable attorneys' fees as available under 17 U.S.C. § 505 or other applicable  
     15 statute.
- 16 e) Destruction or other reasonable disposition of all copies Defendant made or used in  
     17 violation of the exclusive rights of Plaintiffs and the Class, under 17 U.S.C.  
     18 § 503(b).
- 19 f) Pre- and post-judgment interest on the damages awarded to Plaintiffs and the Class,  
     20 and that such interest be awarded at the highest legal rate from and after the date this  
     21 class action complaint is first served on Defendant.
- 22 g) Defendant to pay for the costs and expenses of a Court-approved notice program  
     23 through post and media designed to give immediate notification to the Class.
- 24 h) Further relief for Plaintiffs and the Class as may be just and proper.

25

26

27

28



1 Rohit D. Nath (SBN 316062)  
2 **SUSMAN GODFREY L.L.P**  
3 1900 Avenue of the Stars, Suite 1400  
4 Los Angeles, CA 90067-2906  
5 Telephone: (310) 789-3100  
RNath@susmangodfrey.com

6 Elisha Barron (admitted *pro hac vice*)  
7 Craig Smyser (admitted *pro hac vice*)  
8 **SUSMAN GODFREY L.L.P.**  
9 One Manhattan West, 51st Floor  
10 New York, NY 10019  
11 Telephone: (212) 336-8330  
12 ebarron@susmangodfrey.com  
13 csmysyer@susmangodfrey.com

14 Jordan W. Connors (admitted *pro hac vice*)  
15 Trevor D. Nystrom (admitted *pro hac vice*)  
16 Dylan B. Salzman (admitted *pro hac vice*)  
17 **SUSMAN GODFREY L.L.P**  
18 401 Union Street, Suite 3000  
19 Seattle, WA 98101  
20 Telephone: (206) 516-3880  
jconnors@susmangodfrey.com  
tnystrom@susmangodfrey.com  
dsalzman@susmangodfrey.com

21 Rachel J. Geman (*pro hac vice*)  
22 Danna Z. Elmasry (*pro hac vice*)  
23 **LIEFF CABRASER HEIMANN**  
24 **& BERNSTEIN, LLP**  
25 250 Hudson Street, 8th Floor  
26 New York, NY 10013  
27 Tel.: 212.355.9500  
rgeman@lchb.com  
delmasry@lchb.com

28 Anne B. Shaver  
29 **LIEFF CABRASER HEIMANN**  
30 **& BERNSTEIN, LLP**  
31 275 Battery Street, 29th Floor  
32 San Francisco, CA 94111  
33 Tel.: 415.956.1000  
ashaver@lchb.com

34 Betsy A. Sugar (*pro hac vice*)

**LIEFF CABRASER HEIMANN  
& BERNSTEIN, LLP**  
222 2nd Avenue S. Suite 1640  
Nashville, TN 37201  
Tel.: 615.313.9000  
[bsugar@lchb.com](mailto:bsugar@lchb.com)

David A. Straite (admitted *pro hac vice*)  
**DICELLO LEVITT LLP**  
485 Lexington Avenue, Suite 1001  
New York, NY 10017  
Tel. (646) 933-1000  
[dstraite@dicellosevitt.com](mailto:dstraite@dicellosevitt.com)

Amy E. Keller (admitted *pro hac vice*)  
Nada Djordjevic (admitted *pro hac vice*)  
James A. Ulwick (admitted *pro hac vice*)  
**DiCELLO LEVITT LLP**  
Ten North Dearborn Street, Sixth Floor  
Chicago, Illinois 60602  
Tel. (312) 214-7900  
[akeller@dicellosevitt.com](mailto:akeller@dicellosevitt.com)  
[ndjordjevic@dicellosevitt.com](mailto:ndjordjevic@dicellosevitt.com)  
[julwick@dicellosevitt.com](mailto:julwick@dicellosevitt.com)

Brian O'Mara (SBN 229737)  
**DICELLO LEVITT LLP**  
4747 Executive Drive  
San Diego, California 92121  
Telephone: (619) 923-3939  
Facsimile: (619) 923-4233  
[briano@dicelolevitt.com](mailto:briano@dicelolevitt.com)

Matthew Butterick (State Bar No. 250953)  
1920 Hillhurst Avenue, #406  
Los Angeles, CA 90027  
Telephone: (323) 968-2632  
Facsimile: (415) 395-9940  
[mb@buttericklaw.com](mailto:mb@buttericklaw.com)

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