

1 ~~Joseph R. Saveri (State Bar No. 130064)~~
 2 ~~JOSEPH SAVERI LAW FIRM, LLP~~
 3 ~~601 California Street, Suite 1000~~
 4 ~~San Francisco, CA 94108~~
 5 ~~Telephone: (415) 500-6800~~
 6 ~~Faxsimile: (415) 395-9940~~
 7 ~~Email: jsaveri@saverilawfirm.com~~

8 ~~Matthew Butterick (State Bar No. 250953)~~
 9 ~~1920 Hillhurst Avenue, #406~~
 10 ~~Los Angeles, CA 90027~~
 11 ~~Telephone: (323) 968-2632~~
 12 ~~Faxsimile: (415) 395-9940~~
 13 ~~Email: mb@buttericklaw.com~~

14 ~~Laura M. Matson (pro hac vice pending)~~
 15 ~~LOCKRIDGE GRINDAL NAUEN PLLP~~
 16 ~~100 Washington Avenue South, Suite 2200~~
 17 ~~Minneapolis, MN 55401~~
 18 ~~Telephone: (612) 339-6900~~
 19 ~~Faxsimile: (612) 339-0981~~
 20 ~~Email: lmatson@locklaw.com~~

21 *Counsel for Individual and Representative*
Plaintiffs and the Proposed Class [Counsel on
Signature Page]

22
 23
 24
 25
 26
 27
 28

UNITED STATES DISTRICT COURT
NORTHERN DISTRICT OF CALIFORNIA
SAN FRANCISCO/OAKLAND DIVISION

~~Abdi Nazemian, an individual;~~
~~Brian Keene, an individual; and~~
~~Stewart O'Nan, an individual;~~
~~Andre Dubus III, an individual; and~~
~~Susan Orlean, an individual.~~

Individual and Representative Plaintiffs,

v.

<p>Case No. <u>Master File Case No. 4:24-cv-01454-JST (SK)</u> <u>Consolidated with Case No. 4:24-cv-02655-JST (SK)</u></p> <p>Complaint</p> <p>Class Action</p> <p>Demand for Jury Trial</p>

NVIDIA Corporation, a Delaware corporation;

Defendant.

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

5

6 OVERVIEW

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

9 2. A *large language model* is an AI software program designed to emit convincingly
 10 naturalistic text outputs in response to user prompts. NeMo Megatron-GPT (“NeMo Megatron”) is a
 11 series of large language models created by NVIDIA and released in September 2022.

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

20 4. Plaintiffs and Class members are authors. They own registered copyrights in certain
 21 books that were included in the training dataset that NVIDIA has admitted copying, storing, and using
 22 to traindevelop its NeMo MegatronAI language models. Plaintiffs and Class members never authorized
 23 NVIDIA to use their copyrighted works as training material.

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

27 I am on the data strategy team at NVIDIA, we are exploring including Anna’s
 28 Archive in pre-training data for our LLMs.
 We are figuring out internally whether we are willing to accept the risk of using
 this data, but would like to speak with your team to get a better understanding of
 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.

5.7. 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

6.8. 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).

7.9. Jurisdiction and venue are proper in this judicial district under 28 U.S.C. § 1391(c)(2) 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.

8.10. 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

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

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

11.12. 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*.

13. A nonexhaustive 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*.

14. 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*.

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

DEFENDANT

13.16. 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

14.17. 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.

15.18. 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

1 furtherance thereof. Each acted as the principal, agent, or joint venture of, or for Defendant with respect
 2 to the acts, violations, and common course of conduct alleged herein.

3

4 FACTUAL ALLEGATIONS

5 16.19. NVIDIA is a diversified technology company founded in 1993 that originally focused
 6 on computer-graphics hardware, e.g., Graphics Processing Units (“GPUs”), and has since expanded to
 7 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
 9 programs. NVIDIA’s hardware and software is used by all Frontier AI companies—companies that
 10 develop the most advanced AI systems— which has resulted in NVIDIA becoming the world’s most
 11 valuable company.

12 17.20. In September 2022addition to the hardware and software products it sells to AI
 13 companies, NVIDIA released its NeMo Megatron series ofitself has developed numerous AI models
 14 known as “large language models.A large language model (“” (“LLMs”). An LLM”) is AI software
 15 designed to emit convincingly naturalistic text outputs in response to user prompts. NVIDIA sells
 16 products to its customers that rely on NVIDIA’s LLMs.

17 18.21. Though an LLM is a LLMs are software program, it is programs, they are not created the
 18 way most software programs are—that is, by human software programmers writing code. Rather, an
 19 LLM is LLMs are trained by copying an enormous quantity of textual works and then feeding these
 20 copies in pieces into the model. This corpus of input material is called the *training dataset*.

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

26 19.23. During the training process, LLMs copy and ingest each textual work in the training
 27 dataset and extractextract protected expression from it. TheIn a process somewhat resembling a guess-
 28 and-check quiz, the LLM is progressively adjusts its outputadjusted to more closely approximate the

1 protected expression copied from the training dataset. The LLM records the results of this process in a
 2 large set of numbers called *weights* or *parameters* that are stored within the model, and, in some sense,
 3 “are” the model. These weights are entirely and uniquely derived from the protected expression in the
 4 training dataset. For instance, the NeMo Megatron–GPT 20B language modelmodel—an LLM
 5 released in September 2022 as part of NVIDIA’s NeMo Megatron series of LLMs—is so named
 6 because the model stores 20 billion (“20B”) weights derived from protected expression in its training
 7 dataset.

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

16 20.25. Once the LLM has copied and ingested the textual works in the training dataset and
 17 transformed the protected expression into stored weights, the LLM is able to emit convincing
 18 simulations of natural written language in response to user prompts. Whenever an LLM generates text
 19 output in response to a user prompt, it is performing a computation that relies on these stored weights,
 20 with the goal of imitating the protected expression ingested from the training dataset.

21 21.26. Much of the material in NVIDIA’s training dataset, however, comes from copyrighted
 22 works—including books written by Plaintiffs and Class members—that were acquired, copied and
 23 stored by NVIDIA without consent, without credit, and without compensation.

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

4 22.28. In September 2022, NVIDIA ~~first~~ announced the availability of the NeMo Megatron
 5 language models in a video on its website: “For the first time, NVIDIA is making its checkpoints
 6 available publicly, where the checkpoints are trained with NeMo Megatron … this is just to begin with.
 7 And this is not the end. We will continue to add more checkpoints in the future.”³ In this context
 8 “checkpoints” is an alternate term for language models ~~within the NeMo Megatron series.~~ The
 9 language models released in September 2022 include NeMo Megatron-GPT 1.3B, NeMo Megatron-
 10 GPT 5B, NeMo Megatron-GPT 20B, and NeMo Megatron-T5 3B models.

11 23.29. Each of ~~the~~these NeMo Megatron models ~~is~~was hosted on a website called Hugging
 12 Face, where ~~it has~~a model card ~~that~~provides information about ~~the~~each model, including its training
 13 dataset. The model card for each of the NeMo Megatron models states ~~that~~, “The model was trained on
 14 ‘The Pile’ dataset prepared by EleutherAI.”⁴

15 24.30. The Pile is a training dataset curated by a research organization called EleutherAI. In
 16 December 2020, EleutherAI introduced this dataset in a paper called “The Pile: An 800GB Dataset of
 17 Diverse Text for Language Modeling”⁵ (the “EleutherAI Paper”).

18 25.31. According to the EleutherAI Paper, one of the components of The Pile is a collection of
 19 books called Books3. The EleutherAI Paper reveals that the Books3 dataset comprises 108 gigabytes of
 20 data, or approximately 12% of the dataset, making it the third largest component of The Pile by size.

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

23 ² *Id.*

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

26 ⁴ See, e.g., <https://huggingface.co/nvidia/nemo-megatron-gpt-1.3B#training-data>,
<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 ⁵ Available at <https://arxiv.org/pdf/2101.00027.pdf>

1 26.32. The EleutherAI Paper further describes the contents of Books3:

2
3 Books3 is a dataset of books derived from a copy of the contents of the
4 Bibliotik private tracker ... Bibliotik consists of a mix of fiction and
5 nonfiction books and is almost an order of magnitude larger than our
6 next largest book dataset (BookCorpus2). **We included Bibliotik**
7 **because books are invaluable for long-range context modeling**
8 **research and coherent storytelling.**⁶

9 27.33. Bibliotik is one of a number of notorious “shadow library” websites ~~that also includes~~
10 ~~Library Genesis (aka LibGen), Z Library (aka Book), Sci-Hub which make, store~~, and ~~Anna's Archive~~.
11 ~~These shadow libraries have long been of interest to the AI training community because they host and~~
12 ~~distribute vast huge quantities of unlicensed pirated copyrighted material. For that reason, these~~
13 ~~shadow libraries also violate the U.S. Copyright Act works via the BitTorrent Protocol.~~

14 28.34. The person who assembled the Books3 dataset, Shawn Presser, has confirmed in public
15 statements that it represents “all of Bibliotik” and contains approximately 196,640 books.

16 29.35. Plaintiffs’ copyrighted books listed in Exhibit A are among the works in the Books3
17 dataset. Below, these books are referred to as the **Infringed Works**.

18 30.36. Until October 2023, the Books3 dataset was available from Hugging Face. At that time,
19 the Books3 dataset was removed with a message that it “is defunct and no longer accessible due to
20 reported copyright infringement.”⁷

21 31.37. ~~In sum~~, NVIDIA has publicly admitted training its NeMo Megatron models on a copy
22 of The Pile dataset. Therefore, NVIDIA necessarily also ~~trained its NeMo Megatron models on (1)~~
23 acquired a copy of Books3—~~(because Books3 is part of The Pile)~~ and (2) made additional copies of
24 Books3 during the course of developing LLMs, including (but not limited to) its NeMo Megatron
25 models. Certain books written by Plaintiffs are part of Books3—including the Infringed Works—and
26 thus NVIDIA necessarily ~~trained (1) made unlawful copies of Plaintiffs' works when downloading~~

27 ⁶ *Id.* at 3–4. (emphasis added).

28 ⁷ See https://huggingface.co/datasets/the_pile_books3

1 Books3, and (2) made additional unlawful copies of Plaintiffs' works when developing its LLMs,
 2 including (but not limited to) its NeMo Megatron models~~on one or more copies of the Infringed~~
 3 Works, thereby. NVIDIA thus directly infringing the infringed Plaintiffs' copyrights~~of the Plaintiffs.~~

4 38. But NVIDIA's use of Plaintiffs' Infringed Works was not limited to the models it
 5 publicly disclosed were trained on The Pile. NVIDIA and its engineers maintained The Pile in
 6 centralized servers and repeatedly (and extensively) used The Pile following its acquisition, including to
 7 develop multiple LLMs known internally as NeMo Megatron GPT 126M, NeMo Megatron GPT 40B,
 8 NeMo Megatron GPT 175B, NeMo Megatron T5 220M, NeMo Megatron T5 11B, and NeMo
 9 Megatron T5 23B.

10 39. NVIDIA's use of The Pile to develop language models was not limited to a single line
 11 or class of models either. Instead, language models across NVIDIA used The Pile.

12 40. NVIDIA used The Pile to train and develop models that do not bear the NeMo
 13 Megatron name as well. For instance, NVIDIA included the Pile dataset as training data for an LLM
 14 known as Megatron 345M, which was publicly released as the Megatron GPT2 345m model. NVIDIA
 15 also used The Pile to train an LLM known as "NeMo GPT-3 10B." NVIDIA additionally developed the
 16 InstructRetro-48B and Retro-48B LLMs using the Books3 dataset from The Pile.

17 41. The Pile was not NVIDIA's only dataset that included Books3. NVIDIA also
 18 downloaded the SlimPajama dataset.⁸ "SlimPajama was created by cleaning and deduplicating the 1.2T
 19 token RedPajama dataset from [the company] Together [AI]." And the RedPajama dataset itself
 20 originally included the Books3 dataset. The SlimPajama dataset included the Books3 dataset. NVIDIA
 21 used the SlimPajama dataset to test "both sentencepiece and BPE [tokenizers]." Tokenizers are software
 22 which is used to process training data for use in LLM training and development. In short, NVIDIA used
 23 the SlimPajama dataset to develop and test the software used in the development of its LLMs. As one
 24 NVIDIA employee remarked, "SlimPajama . . . is available in our org." NVIDIA, therefore, again
 25 infringed Plaintiffs' copyrights by downloading unauthorized copies of their works by downloading,
 26 storing, and using the SlimPajama dataset.

27
 28 ⁸ See <https://huggingface.co/datasets/cerebras/SlimPajama-627B>.

1 42. Upon information and belief, NVIDIA also developed a large number of internal
 2 models, including checkpoints, many of which were never given proper names or publications but
 3 which also unlawfully included datasets containing Plaintiffs' and Class members' works, such as The
 4 Pile.

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

7 44. Not content to acquire, store, and use The Pile in its internal and external LLM
 8 research, development, and commercialization efforts, NVIDIA sought vastly more copyrighted works
 9 than The Pile could provide. Because the quality of an LLM depends on both the quality *and* quantity
 10 of its training data, NVIDIA found itself desperate for additional books. Books have the unique
 11 designation of being widely understood as high-quality LLM training data and being available illegally
 12 in large quantities from illicit shadow libraries.

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

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

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

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

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

10 50. As revealed publicly over the last year,⁹ it is an industry-wide practice to use shadow
 11 libraries such as Library Genesis, Z-Library, and Pirate Library Mirror. Virtually every one of the major
 12 LLM developers—including OpenAI, Meta, and Anthropic—pirated books from Library Genesis, Z-
 13 Library, Sci-Hub, and/or Pirate Library Mirror. NVIDIA followed this industry-wide practice and
 14 pirated troves of books from shadow libraries.

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

18
 19 Not too long ago, “shadow-libraries” were dying. Sci-Hub, the massive
 20 illegal archive of academic papers, had stopped taking in new works,
 21 due to lawsuits. “Z-Library”, the largest illegal library of books, saw its
 22 alleged creators arrested on criminal copyright charges . . . *Then came*
 23 *AI. Virtually all major companies building LLMs contacted us to train*
 24 *on our data. . . We have given high-speed access to about 30*
 25 *companies.* <https://annas-archive.org/blog/ai-copyright.html> (emphasis
 26 added.)

26 27 ⁹ See, e.g., Alex Reisner, *The Unbelievable Scale of AI’s Pirated-Books Problem*, The Atlantic (March
 28 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 Anthropic’s use of
 LibGen and Pirate Library Mirror to download millions of copyrighted books).

1 52. Internal documents show competitive pressures drove NVIDIA to piracy. In the fall of
 2 2023, NVIDIA faced a rapidly approaching deadline in the form of its annual developer day. In the year
 3 since the launch of the NeMo Megatron series in September 2022, OpenAI had released ChatGPT to
 4 massive success, resulting in a substantial increase in investor attention on AI. In response, NVIDIA
 5 sought to develop and demonstrate cutting edge LLMs at its fall 2023 developer day. In seeking to
 6 acquire data for what it internally called “NextLargeLLM,” “NextLLMLarge” and “Next Generation
 7 LLM” (collectively, “NextLargeLLM”). NVIDIA was “[h]yper [f]ocused on books corpuses.” NVIDIA
 8 knew that “published books under copyright” are “the most valuable” for developing LLMs and
 9 NVIDIA knew that only books were available in sufficient quantities. And NVIDIA needed to achieve 8
 10 trillion tokens for the “NextLargeLLM,” and books provided this means.

11 53. In August 2023, NVIDIA contacted books publishers to obtain fast “access to large
 12 volumes of unique, high-quality datasets” or “ie. books.” But on information and belief, NVIDIA could
 13 not secure this fast access to the huge quantity of books it needed through publishers. As one book
 14 publisher told NVIDIA, it was “not in a position to engage directly just yet but will be in touch.” In
 15 2023, NVIDIA had “chatted with multiple publishers . . . but none [] wanted to enter into data licensing
 16 deals.”

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

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

26 56. Within a week of contacting Anna’s Archive, and days after being warned by Anna’s
 27 Archive of the illegal nature of their collections, NVIDIA management gave “the green light” to proceed
 28 with the piracy. Anna’s Archive offered NVIDIA millions of pirated copyrighted books. Anna’s Archive

1 also offered access to several million books from Internet Archive, which were only normally available
 2 through Internet Archive's digital lending system (a system which was found to be copyright
 3 infringement by the Second Circuit, *see Hachette Book Grp., Inc. v. Internet Archive*, 115 F.4th 163 (2d
 4 Cir. 2024)). Anna's Archive promised NVIDIA access to "a lot of books," totaling roughly 500 terabytes
 5 of data. By downloading Anna's Archive, NVIDIA pirated additional copies of Plaintiff's Infringed
 6 Works.

7 57. On information and belief, in addition to Anna's Archive and The Pile, NVIDIA also
 8 downloaded books hosted or sourced from other shadow libraries, including LibGen, Sci-Hub, and Z-
 9 Library.

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

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

19 60. Upon information and belief, NVIDIA could not obtain the level of books needed for the
 20 Nemotron models without pirating copyrighted books, including Plaintiffs' Infringed Works.

21 61. In sum, NVIDIA has extensively and repeatedly violated the copyrights of Plaintiffs'
 22 Infringed Works including by acquiring these works from pirated sources, storing them, and enabling
 23 its employees to use them for any purpose, and copying them during the LLM training process.

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

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

4 64. 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 Archive.

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

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

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

19 68. Through the NeMo Megatron Framework and BigNLP platforms, NVIDIA provided
 20 customers with “scripts to automatically download and preprocess The Pile dataset which, until
 21 recently, was hosted externally by Eleuther AI.” Meaning, NVIDIA provided tools and resources for its
 22 customers to use the NVIDIA platform to download The Pile, thereby infringing on Plaintiffs’
 23 copyrights. They scripts were developed to help their customers access these pirated datasets more
 24 quickly and easily. NVIDIA employees expressed concern about the “[t]ime needed for downloading
 25 pile files,” so they developed and distributed code to “download and extract[] 30 pile files [in] ~70
 26 minutes[,] which clearly shows the need for data prep parallelism.”

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

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

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

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

14

COUNT 1

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

15 32.72. Plaintiffs incorporate by reference the preceding factual allegations.

16 33.73. As the owners of the registered copyrights in the Infringed Works, Plaintiffs hold the
 17 exclusive rights to those books under 17 U.S.C. § 106.

18 34.74. To ~~train the NeMo Megatron language models~~ develop NVIDIA's LLMs, NVIDIA
 19 ~~downloaded and~~ copied The Pile ~~dataset~~ and ~~SlimPajama~~ datasets. The Pile ~~dataset includes~~ and
 20 ~~SlimPajama datasets include~~ the Books3 dataset, which includes the Infringed Works. NVIDIA made
 21 multiple copies of the Books3 dataset while ~~training the NeMo Megatron models~~ developing its LLMs.

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

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

35.77. 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.

36.78. 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 the NeMo Megatron 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 for training other models.

37.79. 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.

80. 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

81. Plaintiffs incorporate by reference the preceding factual allegations.

82. 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; controlling or managing the property or other assets with which the direct infringement was accomplished; or providing

1 business, legal, strategic, or operational guidance that allowed its customers to download, copy, and
 2 store Plaintiffs' and Class members' copyrighted works.

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

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

10

11 **COUNT III**

12 **Vicarious Copyright Infringement**

13 **against NVIDIA**

14. Plaintiffs incorporate by reference the preceding factual allegations.

15. 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, to
 18 download The Pile dataset (and potentially other datasets containing copyrighted books as well).
 19 NVIDIA failed to exert its right and ability to control its customers infringing acts.

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

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

CLASS ALLEGATIONS

38.89. The “Class Period” as defined in this Complaint begins ~~on at least~~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.

39.90. 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 domiciled in the United States that own a registered United States copyright in any literary work that was downloaded or otherwise copied by Defendant and / or used as by Defendant in LLM training data for the NeMo Megatron large language models, research, or development during the Class Period.

40.91. 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.

41.92. **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 **42.93. 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 **43.94. Adequacy.** Plaintiffs will fairly and adequately represent the interests of the members
 5 of the Class because the Plaintiffs have experienced the same harms as the members of the Class and
 6 have no conflicts with any other members of the Class. Furthermore, Plaintiffs have retained
 7 sophisticated and competent counsel who are experienced in prosecuting federal and state class actions,
 8 as well as other complex litigation.

9 **44.95. Commonality and predominance.** Numerous questions of law or fact common to each
 10 Class member arise from Defendant's conduct and predominate over any questions affecting the
 11 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 ~~and used them to train the NeMo~~
 14 ~~Megatron language models.~~

15 b. Whether Defendant violated the copyrights of Plaintiffs and the Class when they used
 16 them to research, develop, and train language models.

17 b.c. Whether Defendant intended to cause further infringement of the Infringed Works with
 18 ~~the NeMo Megatron~~these language models because they have distributed these models
 19 under an open license and advertised those models as a base from which to build further
 20 models.

21 d. Whether Defendant's support, facilitation, and encouragement of the infringement by
 22 NVIDIA's customers of Plaintiffs' and Proposed Class Members' copyrighted works
 23 constitutes vicarious or contributory infringement under the Copyright Act

24 e.e. Whether any affirmative defense excuses Defendant's conduct.

25 d.f. Whether any statutes of limitation constrain the potential for recovery for Plaintiffs and
 26 the Class.

27 **45.96. Other class considerations.** Defendant has acted on grounds generally applicable to
 28 the Class. This class action is superior to alternatives, if any, for the fair and efficient adjudication of

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

5

6 **DEMAND FOR JUDGMENT**

7 Wherefore, Plaintiffs request that the Court enter judgment on their behalf and on behalf of the
 8 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 the
 13 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. § 503(b).
- 18 f) Pre- and post-judgment interest on the damages awarded to Plaintiffs and the Class, and
 19 that such interest be awarded at the highest legal rate from and after the date this class
 20 action complaint is first served on Defendant.
- 21 g) Defendant is to be financially responsible to pay for the costs and expenses of a Court-
 22 approved notice program through post and media designed to give immediate
 23 notification to the Class.
- 24 h) Further relief for Plaintiffs and the Class as may be just and proper.

25

26 **JURY TRIAL DEMANDED**

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

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

1 Dated: March 8, 2024

2 By: /s/ Joseph R. Saveri

3 Joseph R. Saveri

4 Joseph R. Saveri (State Bar No. 130064)
5 Christopher K. L. Young (State Bar No. 318371)
6 Elissa Buchanan (State Bar No. 249996)
7 **JOSEPH SAVERI LAW FIRM, LLP**
8 601 California Street, Suite 1000
9 San Francisco, CA 94108
10 Telephone: (415) 500-6800
11 Facsimile: (415) 395-9940
12 Email: jsaveri@saverilawfirm.com
13 cyoung@saverilawfirm.com
14 eabuchanan@saverilawfirm.com

15 Matthew Butterick (State Bar No. 250953)
16 1920 Hillhurst Avenue, #406
17 Los Angeles, CA 90027
18 Telephone: (323) 968-2632
19 Facsimile: (415) 395-9940
20 Email: mb@buttericklaw.com

21 Brian D. Clark (*pro hac vice pending*)
22 Laura M. Matson (*pro hac vice pending*)
23 Arielle S. Wagner (*pro hac vice pending*)
24 Eura Chang (*pro hac vice pending*)
25 **LOCKRIDGE GRINDAL NAUEN PLLP**
26 100 Washington Avenue South, Suite 2200
27 Minneapolis, MN 55401
28 Telephone: (612) 339-6900
Facsimile: (612) 339-0981
Email: bdclark@locklaw.com
lmmatson@locklaw.com
aswagner@locklaw.com
echang@locklaw.com

29 Counsel for Individual and Representative
30 Plaintiffs and the Proposed Class

31 Dated: October 17, 2025

32 By: /s/ Rohit D. Nath

33 Joseph R. Saveri (CSB No. 130064)
34 Christopher K.L. Young (CSB No. 318371)
35 Evan Creutz (CSB No. 349728)
36 Elissa A. Buchanan (CSB No. 249996)
37 William Waldir Castillo Guardado (CSB No. 294159)

1 **JOSEPH SAVERI LAW FIRM, LLP**

2 601 California Street, Suite 1505

3 San Francisco, CA 94108

4 Telephone: (415) 500-6800

5 Facsimile: (415) 395-9940

6 Email: jsaveri@saverilawfirm.com

7 cyoung@saverilawfirm.com

8 ecreutz@saverilawfirm.com

9 eabuchanan@saverilawfirm.com

10 wcastillo@saverilawfirm.com

11 Bryan L. Clobes (admitted pro hac vice)

12 Mohammed Rathur (admitted pro hac vice)

13 **CAFFERTY CLOBES MERIWETHER**

14 **& SPRENGEL LLP**

15 135 South LaSalle Street, Suite 3210

16 Chicago, IL 60603

17 Tel: 312-782-4880

18 bclobes@caffertyclobes.com

19 mrathur@caffertyclobes.com

20 Justin A. Nelson (admitted *pro hac vice*)

21 Alejandra C. Salinas (admitted *pro hac vice*)

22 **SUSMAN GODFREY L.L.P.**

23 1000 Louisiana Street, Suite 5100

24 Houston, TX 77002

25 Telephone: (713) 651-9366

26 Facsimile: (713) 654-6666

27 jnelson@susmangodfrey.com

28 asalinas@susmangodfrey.com

19 Rohit D. Nath (SBN 316062)

20 **SUSMAN GODFREY L.L.P**

21 1900 Avenue of the Stars, Suite 1400

22 Los Angeles, CA 90067-2906

23 Telephone: (310) 789-3100

24 RNath@susmangodfrey.com

25 Elisha Barron (admitted *pro hac vice*)

26 Craig Smyser (admitted *pro hac vice*)

27 **SUSMAN GODFREY L.L.P.**

28 One Manhattan West, 51st Floor

1 New York, NY 10019

2 Telephone: (212) 336-8330

3 ebarron@susmangodfrey.com

4 csmysyer@susmangodfrey.com

5 Jordan W. Connors (admitted *pro hac vice*)

1 Trevor D. Nystrom (admitted *pro hac vice*)
2 SUSMAN GODFREY L.L.P
3 401 Union Street, Suite 3000
4 Seattle, WA 98101
5 Telephone: (206) 516-3880
6 jconnors@susmangodfrey.com
7 tnystrom@susmangodfrey.com

8 Rachel J. Geman (*pro hac vice*)
9 Danna Z. Elmasry (*pro hac vice*)
10 **LIEFF CABRASER HEIMANN**
11 **& BERNSTEIN, LLP**
12 250 Hudson Street, 8th Floor
13 New York, NY 10013
14 Tel.: 212.355.9500
15 rgeman@lchb.com
delmasry@lchb.com

16 Anne B. Shaver
17 **LIEFF CABRASER HEIMANN**
18 **& BERNSTEIN, LLP**
19 275 Battery Street, 29th Floor
20 San Francisco, CA 94111
21 Tel.: 415.956.1000
22 ashaver@lchb.com

23 Betsy A. Sugar (*pro hac vice*)
24 **LIEFF CABRASER HEIMANN**
25 **& BERNSTEIN, LLP**
26 222 2nd Avenue S. Suite 1640
27 Nashville, TN 37201
28 Tel.: 615.313.9000
bsugar@lchb.com

1 David A. Straite (admitted *pro hac vice*)
2 **DiCELLO LEVITT LLP**
3 485 Lexington Avenue, Suite 1001
4 New York, NY 10017
5 Tel. (646) 933-1000
6 dstraite@dicellevitt.com

7 Amy E. Keller (admitted *pro hac vice*)
8 Nada Djordjevic (admitted *pro hac vice*)
9 James A. Ulwick (admitted *pro hac vice*)
10 **DiCELLO LEVITT LLP**
11 Ten North Dearborn Street, Sixth Floor
12 Chicago, Illinois 60602
13 Tel. (312) 214-7900

1 akeller@dicellosevitt.com
2 ndjordjevic@dicellosevitt.com
3 julwick@dicellosevitt.com

4 [Brian O'Mara \(SBN 229737\)](mailto:Brian O'Mara (SBN 229737))
5 DiCELLO LEVITT LLP
6 4747 Executive Drive
7 San Diego, California 92121
8 [Telephone: \(619\) 923-3939](mailto:Telephone: (619) 923-3939)
9 [Facsimile: \(619\) 923-4233](mailto:Facsimile: (619) 923-4233)
10 briano@dicellosevitt.com

11 [Matthew Butterick \(State Bar No. 250953\)](mailto:Matthew Butterick (State Bar No. 250953))
12 1920 Hillhurst Avenue, #406
13 Los Angeles, CA 90027
14 [Telephone: \(323\) 968-2632](mailto:Telephone: (323) 968-2632)
15 [Facsimile: \(415\) 395-9940](mailto:Facsimile: (415) 395-9940)
16 mb@buttericklaw.com

17 [*Counsel for Individual and Representative Plaintiffs*](#)
18 [*and the Proposed Class*](#)