

1 Bryan L. Clobes (*pro hac vice anticipated*)
2 **CAFFERTY CLOBES MERIWETHER &**
3 **SPRENGEL LLP**
4 205 N. Monroe Street,
5 Media, PA 19063
6 Tel: 215-864-2800
7 bclobes@caffertyclobes.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@dicellolevitt.com

6 Amy E. Keller (*pro hac vice anticipated*)
7 Nada Djordjevic (*pro hac vice anticipated*)
8 James A. Ulwick (*pro hac vice anticipated*)
9 **DICELLO LEVITT LLP**
10 Ten North Dearborn Street, Sixth Floor
11 Chicago, Illinois 60602
12 Tel. (312) 214-7900
13 akeller@dicellolevitt.com
14 ndjordjevic@dicellolevitt.com
15 julwick@dicellolevitt.com

Attorneys for Plaintiffs and the Class

12 David A. Straite (*pro hac vice anticipated*)
13 **DICELLO LEVITT LLP**
14 485 Lexington Avenue, Suite 1001
15 New York, NY 10017
16 Tel. (646) 933-1000
17 dstraite@dicellolevitt.com

17 **UNITED STATES DISTRICT COURT**
18 **NORTHERN DISTRICT OF CALIFORNIA - SAN FRANCISCO DIVISION**

19 ANDRE DUBUS III, SUSAN ORLEAN,

Case No.

20 Individually and on behalf of all others
21 similarly situated,

CLASS ACTION COMPLAINT

22 v.

DEMAND FOR JURY TRIAL

23 NVIDIA CORPORATION, a Delaware
24 corporation;

25 Defendant.
26
27
28

1 Plaintiffs Andre Dubus III and Susan Orlean (“Plaintiffs”), on behalf of themselves and
2 all others similarly situated, bring this class-action complaint (“Complaint”) against defendant
3 NVIDIA Corporation (“NVIDIA” or “Defendant”).
4

5 INTRODUCTION

6 1. *Artificial intelligence*—commonly abbreviated “AI”—denotes software that is
7 designed to algorithmically simulate human reasoning or inference, often using statistical
8 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
11 Megatron”) is a series of large language models created by NVIDIA and released in September
12 2022.

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

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

25 5. NVIDIA copied Plaintiffs’ and Class members’ copyrighted works multiple times
26 to train its NeMo Megatron language models.
27
28

AGENTS AND CO-CONSPIRATORS

1
2 13. The unlawful acts alleged against Defendant in this class action complaint were
3 authorized, ordered, or performed by the Defendant’s respective officers, agents, employees,
4 representatives, or shareholders while actively engaged in the management, direction, or
5 control of the Defendant’s businesses or affairs. The Defendant’s agents operated under the
6 explicit and apparent authority of their principals. Each Defendant, and its subsidiaries,
7 affiliates, and agents operated as a single unified entity.

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

FACTUAL ALLEGATIONS

12
13
14 15. NVIDIA is a diversified technology company founded in 1993 that originally
15 focused on computer-graphics hardware and has since expanded to other computationally
16 intensive fields, including software and hardware for training and operating AI software
17 programs.

18 16. In September 2022, NVIDIA released its NeMo Megatron series of *large language*
19 *models*. A large language model (“LLM”) is AI software designed to emit convincingly naturalistic
20 text outputs in response to user prompts.

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

25 18. During training, the LLM copies and ingests each textual work in the training
26 dataset and extracts protected expression from it. The LLM progressively adjusts its output to
27 more closely approximate the protected expression copied from the training dataset. The LLM
28 records the results of this process in a large set of numbers called *weights* that are stored within the

1 model. These weights are entirely and uniquely derived from the protected expression in the
2 training dataset. For instance, the NeMo Megatron–GPT 20B language model is so named
3 because the model stores 20 billion (“20B”) weights derived from protected expression in its
4 training dataset.

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

10 20. Much of the material in NVIDIA’s training dataset, however, comes from
11 copyrighted works—including books written by Plaintiffs and Class members—that were copied
12 by NVIDIA without consent, without credit, and without compensation.

13 21. In September 2022, NVIDIA first announced the availability of the NeMo
14 Megatron language models in a video on its website: “For the first time, NVIDIA is making its
15 checkpoints available publicly, where the checkpoints are trained with NeMo Megatron ... this is
16 just to begin with. And this is not the end. We will continue to add more checkpoints in the future.”¹
17 In this context “checkpoints” is an alternate term for language models within the NeMo Megatron
18 series. The language models released in September 2022 include NeMo Megatron-GPT 1.3B,
19 NeMo Megatron-GPT 5B, NeMo Megatron-GPT 20B, and NeMo Megatron-T5 3B.

20 22. Each of the NeMo Megatron models is hosted on a website called Hugging Face.
21 Each of the NeMo Megatron models has a *model card* that provides information about the model,
22 including its training dataset. The model card for each of the NeMo Megatron models states that,
23 “The model was trained on ‘The Pile’ dataset prepared by EleutherAI.”²

24
25
26 ¹ See <https://www.nvidia.com/en-us/on-demand/session/gtcfall22-a41200/?nvid=nv-int-tblg-881125>,
starting at 37:25.

27 ² See, e.g., <https://huggingface.co/nvidia/nemo-megatron-gpt-1.3B#training-data>,
28 <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>

1 23. The Pile is a training dataset curated by a research organization called EleutherAI. In
2 December 2020, EleutherAI introduced this dataset in a paper called “The Pile: An 800GB Dataset
3 of Diverse Text for Language Modeling”³ (the “EleutherAI Paper”).

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

8 25. The EleutherAI Paper further describes the contents of Books3:

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

13 26. Bibliotik is one of a number of notorious “shadow library” websites that also
14 includes Library Genesis (aka LibGen), Z-Library (aka B-ok), Sci-Hub, and Anna’s Archive.
15 These shadow libraries have long been of interest to the AI-training community because they
16 host and distribute vast quantities of unlicensed copyrighted material. For that reason, these
17 shadow libraries also violate the U.S. Copyright Act.

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

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

22 29. Until October 2023, the Books3 dataset was available from Hugging Face. At that
23 time, the Books3 dataset was removed with a message that it “is defunct and no longer
24 accessible due to reported copyright infringement.”⁵
25

26
27 ³ Available at <https://arxiv.org/pdf/2101.00027.pdf>

28 ⁴ *Id.* at 3-4.

⁵ See https://huggingface.co/datasets/the_pile_books3

1 States. Therefore, joinder of all members of the Class in the prosecution of this action is
2 impracticable.

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

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

12 44. **Commonality and predominance.** Numerous questions of law or fact common
13 to each Class member arise from Defendant's conduct and predominate over any questions
14 affecting the members of the Class individually:

- 15 a. Whether Defendant violated the copyrights of Plaintiffs and the Class when they
16 obtained copies of Plaintiffs' and Class members' Infringed Works and copied the
17 Infringed Works into the dataset used to train the NeMo Megatron language
18 models.
- 19 b. Whether Defendant intended to cause further infringement of the Infringed Works
20 with the NeMo Megatron language models because they have distributed these
21 models under an open license and advertised those models as a base from which
22 to build further models.
- 23 c. Whether any affirmative defense excuses Defendant's conduct.
- 24 d. Whether any statutes of limitation limits the potential for recovery for Plaintiffs
25 and the Class.

26 45. **Other class considerations.** Defendant has acted on grounds generally
27 applicable to the Class. This class action is superior to alternatives, if any, for the fair and
28 efficient adjudication of this controversy. Prosecuting the claims pleaded herein as a class action

1 will eliminate the possibility of repetitive litigation. There will be no material difficulty in the
2 management of this action as a class action. The prosecution of separate actions by individual
3 Class members would create the risk of inconsistent or varying adjudications, establishing
4 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
8 of 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
17 in 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
20 Class, and that such interest be awarded at the highest legal rate from and after
21 the date this class action complaint is first served on Defendant.
- 22 g) Defendant is to be jointly and severally responsible financially for the costs and
23 expenses of a Court-approved notice program through post and media designed
24 to give immediate notification to the Class.
- 25 h) Further relief for Plaintiffs and the Class as may be just and proper.
- 26
27
28

JURY TRIAL DEMANDED

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

Dated: May 2, 2024

By: /s/ Brian O'Mara
BRIAN O'MARA, SBN 229737
briano@dicellolevitt.com
DiCELLO LEVITT LLP
4747 Executive Drive
San Diego, California 92121
Telephone: (619) 923-3939
Facsimile: (619) 923-4233

Bryan L. Clobes (*pro hac vice anticipated*)
**CAFFERTY CLOBES MERIWETHER
& SPRENGEL LLP**
205 N. Monroe Street
Media, PA 19063
Tel: 215-864-2800
bclobes@caffertyclobes.com

Alexander J. Sweatman (*pro hac vice anticipated*)
**CAFFERTY CLOBES MERIWETHER
& SPRENGEL LLP**
135 South LaSalle Street, Suite 3210
Chicago, IL 60603
Tel: 312-782-4880
asweatman@caffertyclobes.com

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

David A. Straite (*pro hac vice anticipated*)
DiCELLO LEVITT LLP
485 Lexington Avenue, Suite 1001

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

New York, NY 10017
Tel. (646) 933-1000
dstraite@dicellolevitt.com

*Attorneys for Plaintiffs and the Proposed
Class*

EXHIBIT A

The Garden of Last Days (TX0006864182)

Type of Work: Text

Registration Number / Date:
TX0006864182 / 2008-07-16

Application Title: The Garden of Last Days: A Novel.

Title: The Garden of Last Days: A Novel.

Description: Book, 537 p.

Copyright Claimant:
Andre Dubus III.

Date of Creation: 2007

Date of Publication:
2008-06-02

Nation of First Publication:
United States

Authorship on Application:
Andre Dubus III; Domicile: United States. Authorship:
text.

Copyright Note: Regarding material excluded: Deposit states reprint
material used with permission

ISBN: 978-0-393-04165-1

Names: Dubus III, Andre

The Cage Keeper (TX0002495806)

Type of Work: Text

Registration Number / Date:
TX0002495806 / 1989-02-14

Title: The Cage keeper and other stories / Andre Dubus III.

Edition: 1st ed.

Imprint: New York : Dutton, c1989.

Description: 184 p.

Copyright Claimant:
Andre Dubus III

Date of Creation: 1989

Date of Publication:
1989-01-30

Variant title: The Cage keeper

Names: Dubus, Andre 3rd, 1959-

Townie: A Memoir (TX0007344763)

Type of Work: Text

Registration Number / Date:
TX0007344763 / 2011-04-15

Application Title: Townie: A Memoir.

Title: Townie: A Memoir.

Description: Book, 387 p.

Copyright Claimant:
Andre Dubus III.

Date of Creation: 2010

Date of Publication:
2011-02-23

Nation of First Publication:
United States

Authorship on Application:
Andre Dubus III; Domicile: United States. Authorship:
text,
editing.

Pre-existing Material:
text.

Basis of Claim: text, editing.

ISBN: 978-0-393-06466-7

Names: Dubus, Andre, III

The Orchid Thief (TX0004921990)

Type of Work: Text

Registration Number / Date:
TX0004921990 / 1999-02-24

Title: The orchid thief / Susan Orlean.

Edition: 1st ed,

Imprint: New York : Random House, c1998.

Description: 282 p.

Copyright Claimant:
Susan Orlean

Date of Creation: 1998

Date of Publication:
1998-12-04

Previous Registration:
Appl. states entire text new except quotes from other
sources.

Names: Orlean, Susan
