

1 Elizabeth J. Cabraser (State Bar No. 083151)  
Daniel M. Hutchinson (State Bar No. 239458)  
2 Reilly T. Stoler (State Bar No. 310761)  
LIEFF CABRASER HEIMANN &  
3 BERNSTEIN, LLP  
275 Battery Street, 29th Floor  
4 San Francisco, CA 94111-3339  
Telephone: (415) 956-1000  
5 ecabraser@lchb.com  
dhutchinson@lchb.com  
6 rstoler@lchb.com

7 Rachel Geman (*pro hac vice forthcoming*)  
LIEFF CABRASER HEIMANN &  
8 BERNSTEIN, LLP  
250 Hudson Street, 8th Floor  
9 New York, New York 10013-1413  
Telephone: (212) 355-9500  
10 rgeman@lchb.com

11 Scott J. Sholder (*pro hac vice forthcoming*)  
CeCe M. Cole (*pro hac vice forthcoming*)  
12 COWAN DEBAETS ABRAHAMS &  
SHEPPARD LLP  
13 60 Broad Street, 30th Floor  
New York, New York 10004  
14 Telephone: (212) 974-7474  
ssholder@cdas.com  
15 ccole@cdas.com

16 *Attorneys for Plaintiff and the Proposed Class*

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

21 CHRISTOPHER FARNSWORTH,  
22 *individually and on behalf of others*  
23 *similarly situated,*

24 Plaintiff,

25 v.

26 META PLATFORMS, INC.,

27 Defendant.  
28

Case No. 3:24-cv-6893

**CLASS ACTION COMPLAINT**

**JURY TRIAL DEMANDED**

1 Plaintiff Christopher Farnsworth, on behalf of himself and all other similarly situated  
2 individuals (the “Class,” as defined below), brings this Class Action Complaint against Defendant  
3 Meta Platforms, Inc. (“Meta”).

4 **NATURE OF THE ACTION**

5 1. Meta stole hundreds of thousands of pirated copyrighted books to build a  
6 commercial product called a Large Language Model (“LLM”). The United States Constitution  
7 recognizes the fundamental principle that copyright holders are entitled to exclusive rights in their  
8 works in order to incentivize further creation. Meta ignored this basic principle, and the federal  
9 law embodying it, by stealing and exploiting copyright-protected books for profit.

10 2. Plaintiff Christopher Farnsworth is a best-selling fiction author. He brings this  
11 class action under the Copyright Act to redress the harm Meta’s infringement caused to him and  
12 other authors. Meta exploited Plaintiff’s works without authorization, made illicit copies of them,  
13 and then fed those copies to its LLM. Meta did this to enhance the quality of its LLM’s language  
14 output and, ultimately to have a more desirable and profitable LLM product to compete in the AI  
15 arms race that has developed over the last few years.

16 3. Meta calls its set of LLMs “Llama.” Like other LLMs, Llama is an AI software  
17 program designed to emit convincingly naturalistic text outputs in response to user prompts. The  
18 quality of Llama’s output is determined by the quality of the data it is fed. Meta knew that the  
19 more high-quality, long-form text it fed to Llama, the better Llama would perform commercially.  
20 In this way, Llama is what it ingests. By willfully including unauthorized and pirated copies of  
21 Plaintiff’s copyrighted literary works in its development work, pre-training, and training data,  
22 Meta exploited authors’ literary talents without consent or compensation of any kind.

23 **JURISDICTION AND VENUE**

24 4. The Court has subject matter jurisdiction under 28 U.S.C. §§ 1331 and 1338(a)  
25 because this action arises under the Copyright Act of 1976, 17 U.S.C. § 101, *et seq.*

26 5. The Court also has personal jurisdiction over Defendant because it has purposely  
27 availed itself of the privilege of conducting business in this District.

28



1 LLMs were later renamed Llama 1. Meta initially released Llama 1 for non-commercial research  
2 uses.

3 13. Less than five months later, in July of 2023, Meta released Llama 2, this time  
4 “available under a permissive commercial license.”<sup>1</sup> By September 27, 2023, Meta reported that  
5 users had downloaded over 30 million copies of Llama-based models.<sup>2</sup>

6 14. On April 18, 2024, Meta released Llama 3, which it built on the foundation of  
7 Llama 1 and 2.<sup>3</sup> Upon its commercial release, Meta touted Llama 3 as powering its new consumer  
8 service Meta AI, “one of the world’s leading AI assistants.”<sup>4</sup> Meta admittedly is using its LLMs  
9 to enhance its current commercial products and, on information and belief, Meta is working on a  
10 premium, paid-subscription version of its AI assistant service—powered by Llama 3.

## 11 **II. Meta Developed Its Commercial AI Models Using Stolen, Copyrighted Material.**

### 12 **A. Large Language Models and the Training Process.**

13 15. At issue is a species of AI models called large language models, or LLMs. LLMs  
14 are designed to mimic human use of language. LLMs are able to simulate patterns of human  
15 language by processing input text (“prompts”) and generating output text in response to these  
16 prompts on a predictive basis, *i.e.*, determining what word follows what.

17 16. At a high level, LLMs are algorithms designed to distill mathematically the  
18 relationships between words in written works through a process called “training.” LLMs achieve  
19 this goal by ingesting massive amounts of training materials such as books, breaking down input  
20 text into smaller pieces—words or portions of words, called “tokens”—then translating those  
21 pieces into “vectors,” or a sequence of numbers that is used to identify the token within the series  
22 of algorithms. Those vectors help place each token in a probabilistic context identifying other  
23 tokens closely associated with the word. As described by industry-leading generative AI  
24 development company OpenAI, in comments to the Copyright Office, “the process begins by  
25 breaking text down into roughly word-length ‘tokens,’ which are converted to numbers. The

26 \_\_\_\_\_  
27 <sup>1</sup> <https://ai.meta.com/llama/faq/>

28 <sup>2</sup> <https://ai.meta.com/blog/llama-2-updates-connect-2023/>

<sup>3</sup> <https://arxiv.org/pdf/2407.21783>

<sup>4</sup> <https://about.fb.com/news/2024/04/meta-ai-assistant-built-with-llama-3/>

1 model then calculates each token’s proximity to other tokens in the training data—essentially,  
2 how near one word appears in relation to any other word. These relationships between words  
3 reveal which words have similar meanings . . . and functions.”<sup>5</sup>

4 17. As a model trains by digesting more and more written works, the algorithms,  
5 which distill the relationship between various tokens, changes with it. A model takes text inputs  
6 in the form of an incomplete phrase or passage, and attempts to complete the phrase, essentially a  
7 fill-in-the-blank quiz. A model compares its predicted phrase completion with the actual “correct”  
8 answer and then adjusts its internal algorithms to “learn” from its mistakes and minimize the  
9 difference between any given text input and the “correct” text output.

10 18. A model then repeats this same cycle millions, possibly billions, of times across  
11 the entire corpus of training materials, adjusting its algorithms each time to reflect the text input  
12 from the corpus. This is known as the “pre-training” process, which is the foundation of creating  
13 a “base” LLM model, which can be “fine-tuned” later to achieve more specific results. The pre-  
14 training process fundamentally enables a model to process prompts and generate text output that  
15 mimics human language. It does so by exposing a model to a wide range of texts and using  
16 algorithms to predict the next word in the text. By repeating this process, a model develops  
17 fluency in style, syntax, and expression of ideas, largely by digesting and processing the protected  
18 expression contained in the material used for training. The LLM mines the expression contained  
19 in the training corpus, adjusting its algorithms such that it can mimic the ordering of words, style,  
20 syntax, and presentation of facts, concepts, and themes.

21 19. In a literal sense, a model is what it ingests: without training on material, there is  
22 no LLM. The quality and quantity of the training corpus is critical to the quality of the resulting  
23 model. As one researcher put it: “[large language] model behavior is not determined by  
24 architecture, hyperparameters, or optimizer choices [i.e. technical features set during model  
25 training]. *It’s determined by your dataset, nothing else. Everything else is a means to an end in*  
26

---

27 <sup>5</sup> See Comment of OpenAI “*Re: Notice of Inquiry and Request for Comment* [Docket No. 2023-  
28 06],” United States Copyright Office, Oct. 30, 2023, p. 5-6 (available at:  
[https://downloads.regulations.gov/COLC-2023-0006-8906/attachment\\_1.pdf](https://downloads.regulations.gov/COLC-2023-0006-8906/attachment_1.pdf)).

1 *efficiently deliver[ing] compute to approximating that dataset.*”<sup>6</sup> Conversely, if the LLM’s  
2 training materials are of poor quality, the output—the end product—also will be of poor quality.

3 20. Books are especially valuable training material for the training and development of  
4 LLMs. As one commentator put it, “[b]ooks offer formal and lengthy texts which help LLMs  
5 understand complex language structures, grasp long-term context, and produce coherent  
6 narratives.”<sup>7</sup>

7 21. Any LLM training process, including Meta’s training of Llama, involves creating  
8 multiple copies of the training text which often include a vast number of written works in their  
9 entirety. As the U.S. Patent and Trademark Office has observed, LLM “training” “almost by  
10 definition involve[s] the reproduction of entire works or substantial portions thereof.”<sup>8</sup>

11 **B. Meta Copied A Massive Trove of Pirated Books To Train Its Llama Models.**

12 22. In public disclosures, Meta admitted that it made unauthorized copies and willfully  
13 reproduced copies of nearly two hundred thousand pirated copyrighted books to advance its  
14 commercial AI training and development projects.

15 23. Since those disclosures, Meta has become less forthright about its AI work, but  
16 what Meta previously disclosed confirms that Meta downloaded and copied a dataset of text  
17 without authorization of registered copyright owners, including Plaintiff, called The Pile as part  
18 of its work in training and developing its LLMs.

19 24. The Pile is an 800 GB open-source dataset created for training large language  
20 models that is, and was well known by AI developers. At the time Meta downloaded The Pile, it  
21 was hosted and made publicly available online by a nonprofit called EleutherAI. As described by  
22 its creators, “The Pile is constructed from 22 diverse high-quality subsets . . . many of which  
23

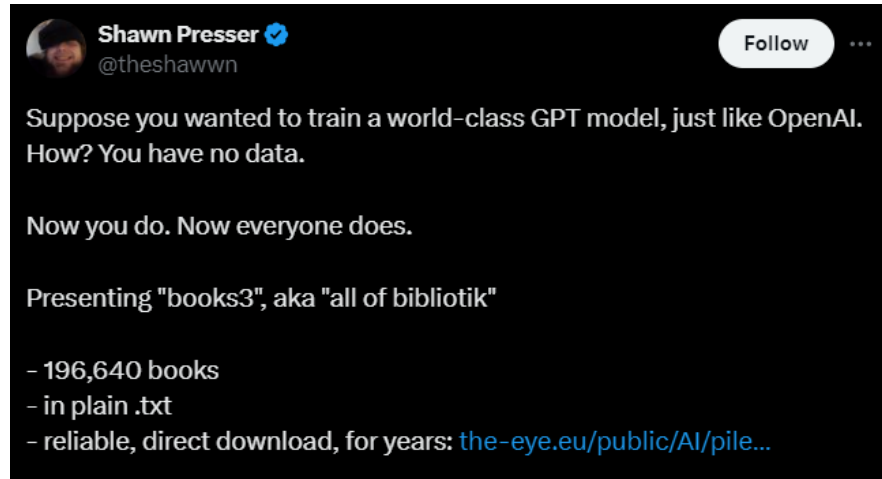
24 <sup>6</sup> See “The ‘it’ in AI models is the data set,” James Betker, <https://nonint.com/2023/06/10/the-it-in-ai-models-is-the-dataset/> (June 10, 2023) (emphasis added).

25 <sup>7</sup> “Pretraining of Large Language Models,” Ritwik Raha, <https://gist.github.com/ritwikraha/77e79990992043f60a9588610b2781c5> (last accessed Sept. 25, 2024).

26 <sup>8</sup> U.S. Patent & Trademark Office, *Public Views on Artificial Intelligence and Intellectual*  
27 *Property Policy* 29 (2020), available at  
28 [https://www.uspto.gov/sites/default/files/documents/USPTO\\_AI-Report\\_2020-10-07.pdf](https://www.uspto.gov/sites/default/files/documents/USPTO_AI-Report_2020-10-07.pdf) (last accessed Jan. 22, 2024).

1 derive from academic and professional sources. . . . [M]odels trained on the Pile improve  
 2 significantly over Raw CC and CC-100 on all components of the Pile, which improving  
 3 performance on downstream evaluations.”<sup>9</sup>

4 25. One of The Pile’s architects is an independent developer named Shawn Presser.  
 5 Presser created a dataset included in The Pile called “Books3,” which is a trove of pirated books.  
 6 Presser described how he created Books3 in a Twitter thread from October 2020:<sup>10</sup>



15 26. Presser explained he created Books3 in response to “OpenAI’s papers on GPT-2  
 16 and 3,” which “references to datasets named ‘books1’ and ‘books2,’” the latter of which Presser  
 17 suspects “might be ‘all of libgen.’”<sup>11</sup> LibGen refers to “Library Genesis,” a website offering  
 18 pirated books that was ordered shut down for copyright infringement in 2015. *See Elsevier, Inc. et*  
 19 *al. v. www.Sci-hub.org et al.*, 15-cv-2482-RWS, Dkt. No. 53 (Oct. 30, 2015).

20 27. To create a pirated-book dataset comparable to what he suspected OpenAI created  
 21 for itself, Presser announced that Books3 was a direct download of all “196,640 books” from a  
 22 *different* pirate website called “bibliotik.”<sup>12</sup>

23

24

25 <sup>9</sup> “The Pile: An 800GB Dataset of Diverse Text for Language Modeling,” Gao et al, Abstract,  
<https://arxiv.org/pdf/2101.00027> (last accessed Sept. 25, 2024) (“Gao Paper”).

26 <sup>10</sup> *See* Tweet by Shawn Presser, Oct. 25, 2020,  
<https://x.com/theshawwn/status/1320282149329784833?lang=en> (last accessed Sept. 25, 2024).

27 <sup>11</sup> *See* Tweet by Shawn Presser, Oct. 25, 2020,  
<https://x.com/theshawwn/status/1320282152689336320> (last accessed Sept. 25, 2024).

28 <sup>12</sup> *See* Tweet by Shawn Presser, Oct. 25, 2020,  
<https://x.com/theshawwn/status/1320282149329784833?lang=en> (last accessed Aug. 15, 2024).

1 28. Bibliotik is a “notorious pirated collection” of “pirated books.”<sup>13</sup> For years prior  
2 to its use as “Books3,” Bibliotik was frequently included in roundups of the best—and most  
3 popular—sources for pirated books.<sup>14</sup>

4 29. Books3 was a critical part of The Pile. In EleutherAI’s paper on The Pile, it  
5 explained the key value of Books3 as training material:

6 Books3 is a dataset of books derived from a copy of the contents of the  
7 Bibliotik private tracker . . . Bibliotik consists of a mix of fiction and  
8 nonfiction books and is almost an order of magnitude larger than our next  
9 largest book dataset (BookCorpus2). *We included Bibliotik because books  
are invaluable for long-range context modeling research and coherent  
storytelling.*<sup>15</sup>

10 30. Presser and EleutherAI repeatedly and publicly acknowledged that The Pile and  
11 Books3 was a cache of pirated material. EleutherAI’s paper on The Pile noted that “there is little  
12 acknowledgment of the fact that the processing and distribution of data owned by others may also  
13  
14

---

15 <sup>13</sup> See Schoppert, “Whether you’re an undergraduate doing research, or a fan of the Nick Stone  
16 novels, or indeed a hungry AI...,” Nov. 29, 2022, [https://aicopyright.substack.com/p/whether-  
youre-an-undergraduate-doing](https://aicopyright.substack.com/p/whether-youre-an-undergraduate-doing) (“What is Bibliotik? A notorious pirated collection.”); “What I  
17 Found in a Database Meta Uses to Train Generative AI,” Alex Reisner, *The Atlantic*, Sept. 25,  
2023, [https://www.theatlantic.com/technology/archive/2023/09/books3-ai-training-meta-  
copyright-infringement-lawsuit/675411/](https://www.theatlantic.com/technology/archive/2023/09/books3-ai-training-meta-copyright-infringement-lawsuit/675411/) (“a collection of pirated ebooks, most of them published  
18 in the past 20 years.”); “Revealed: The Authors Whose Pirated Books are Powering Generative  
AI,” Alex Reisner, *The Atlantic*, Aug. 19, 2023,  
19 [https://www.theatlantic.com/technology/archive/2023/08/books3-ai-meta-llama-pirated-  
books/675063/](https://www.theatlantic.com/technology/archive/2023/08/books3-ai-meta-llama-pirated-books/675063/) (“collections of pirated books, such as Library Genesis, Z-Library, and Bibliotik,  
20 that circulate via the BitTorrent file-sharing network.”); “Are ChatGPT, Bard and Dolly 2.0  
Trained On Pirated Content?,” Roger Monti, *Search Engine Journal*, April 20, 2023,  
21 [https://www.searchenginejournal.com/are-chatgpt-bard-and-dolly-2-0-trained-on-pirated-  
content/485089/](https://www.searchenginejournal.com/are-chatgpt-bard-and-dolly-2-0-trained-on-pirated-content/485089/) (“The Books3 dataset contains the text of books that were pirated and hosted at a  
22 pirate site called, bibliotik.”).

23 <sup>14</sup> See Commit History of “Awesome Piracy,” *Github.com*, Oct. 13, 2018,  
[https://github.com/aviranzerioniac/awesome-  
piracy/commit/61928765e3ee0b4f3dbe3c0724b196e5f0f17e59?short\\_path=5a831ea#diff-  
5a831ea67cf5cf8703b0de46901ab25bd191f56b320053be9332d9a3b0d01d15](https://github.com/aviranzerioniac/awesome-piracy/commit/61928765e3ee0b4f3dbe3c0724b196e5f0f17e59?short_path=5a831ea#diff-5a831ea67cf5cf8703b0de46901ab25bd191f56b320053be9332d9a3b0d01d15) (October 13, 2018  
24 commit to “awesome piracy” repo listing “Bibliotik Popular ebooks/audiobooks private  
25 tracker”); “Reddit Piracy Megathread” repo., *Github.com*, Mar. 21, 2019  
[https://github.com/magicoflois/Reddit-Piracy-  
Megathread/blob/master/data/findingtextbooks.md](https://github.com/magicoflois/Reddit-Piracy-Megathread/blob/master/data/findingtextbooks.md) (March 21, 2019 guide from “r/piracy” on  
26 how to source textbooks listing “Bibliotik”); “List of free eBook download sites,” *Pirates-  
forum.org*, Mar. 06 2014, [https://pirates-forum.org/Thread-List-of-free-eBook-download-  
sites?highlight=bibliotik](https://pirates-forum.org/Thread-List-of-free-eBook-download-sites?highlight=bibliotik) (March 31, 2014 post from “pirates forum” thread entitled “List of free  
27 eBook download sites” listing “bibliotik”).  
28

<sup>15</sup> Gao Paper, p. 3-4.

1 be a violation of copyright law.”<sup>16</sup> Furthermore, The Pile’s datasheet notes that “Books3 is almost  
 2 entirely comprised of copyrighted works . . .”<sup>17</sup> Presser, for his part, admitted to releasing  
 3 Books3 despite “fear of copyright backlash”<sup>18</sup>

4 31. In August 2023, Books3 was removed from the “most official” copy of The Pile—  
 5 hosted by “The Eye,” an online repository of training datasets—due to copyright complaints.  
 6 Despite this removal, the original version appears otherwise available as part of The Pile from  
 7 other sources.

8 32. In its February 27, 2023 research paper, Meta admitted to downloading and  
 9 reproducing Books3 willfully as part of its Llama development project. Meta AI researchers  
 10 explained that the training data used to develop Llama 1 included “two book corpora . . . the  
 11 Gutenberg Project, which contains books that are in the public domain, and the Books3 section of  
 12 The Pile (Gao et al., 2020), a publicly available dataset for training large language models.”<sup>19</sup>

13 33. In referring to the Books3 dataset it used, Meta researchers: (1) directly cited the  
 14 2020 EleutherAI paper (authored by Leo Gao) which describes Books3 as “a dataset of books  
 15 derived from a copy of the contents of the Bibliotik private tracker;”<sup>20</sup> and (2) distinguished  
 16 between the public *domain* books in the Gutenberg dataset and the publicly *available* (i.e.  
 17 unauthorized and pirated) books in the Books3 dataset. Meta knew Books3 was a trove of  
 18 copyrighted content sourced from pirate websites like Bibliotik and used it anyway.

19 34. Notably, in an earlier research paper authored by Meta AI in June 2022, titled  
 20 “OPT: Open Pre-trained Transformer Language Models,” Meta indicated that it had eliminated  
 21 many of The Pile component datasets from its LLM pre-training data because “[Meta] found they  
 22 increased the risk of instabilities . . . *or were otherwise deemed unsuitable.*”<sup>21</sup>

23  
 24 \_\_\_\_\_  
 25 <sup>16</sup> *Id.* at 14-15.

26 <sup>17</sup> “Datasheet for the Pile,” Gao et al, Jan. 20, 2022, p. 15, <https://arxiv.org/pdf/2201.07311> (last  
 27 accessed Sept. 25, 2024).

28 <sup>18</sup> Comment of “sillysaurusx,” *Hacker News*, Jul. 11, 2023,  
<https://news.ycombinator.com/item?id=36685115> (last accessed Sept. 25, 2024).

<sup>19</sup> <https://arxiv.org/pdf/2302.13971>

<sup>20</sup> Gao Paper, p. 3

<sup>21</sup> <https://arxiv.org/pdf/2205.01068>

1           35. For the purposes of the work memorialized in this June 2022 research paper, Meta  
2 excised component datasets that included copyrighted material.

3           36. For example, Meta excluded the “PubMed Central” and “PhilPapers” datasets.  
4 PubMed Central is a “repository for biomedical articles run by the United States of America’s  
5 National Center for Biotechnology Information (NCBI).”<sup>22</sup> The PubMed Central dataset contains  
6 copyrighted material.

7           37. Likewise, the PhilPapers dataset is a repository of “*open-access* philosophy  
8 publications from an international database maintained by the Center for Digital Philosophy at the  
9 University of Western Ontario.”<sup>23</sup> The PhilPapers dataset contains copyrighted material.

10           38. In June 2022, Meta also excluded—as “unsuitable”—the Books3 dataset from its  
11 training data. Shortly after that, something changed and as noted above, by February 2023, Meta  
12 had copied Books3 and used it to train its LLMs.

13           39. In March 2023, the Llama 1 language models were leaked to a public internet site  
14 and have continued to circulate.

15           40. Later in March 2023, Meta issued a DMCA takedown notice to a programmer that  
16 released a tool to help users download the leaked Llama 1 language models. In the notice, Meta  
17 asserted copyright over the Llama 1 large language models.

18           41. In July of 2023—after Meta was sued for its infringing use of copyrighted material  
19 for LLM training—Meta released Llama 2, which it described in a July 19, 2023 research paper  
20 entitled “Llama 2: Open Foundation and Fine-Tuned Chat Models”<sup>24</sup> (the Llama 2 Paper).

21           42. The Llama 2 paper does not indicate the specific datasets used for pre-training but  
22 indicates that *more* data was used (not less) and again conspicuously describes the data it sourced  
23 as “publicly available online,” and not public domain.

24           43. Meta’s ongoing exploitation of Books3 was particularly egregious because it took  
25 place after Books3 was taken down from The Pile in 2023 due to copyright complaints directly  
26

---

27 <sup>22</sup> Gao Paper, p. 3

28 <sup>23</sup> Gao Paper, p. 5 (emphasis added).

<sup>24</sup> <https://arxiv.org/pdf/2307.09288.pdf>

1 before Meta sought DMCA removal of the Llama 1 model to vindicate its own copyrights. This is  
2 especially ironic considering Llama 1 is a commercial LLM that was developed—without  
3 authorization—on pirated copies of hundreds of thousands of copyrighted literary works,  
4 including the literary works of Plaintiff.

5 44. Instead of willfully downloading and reproducing a notorious trove of pirated  
6 material, Meta could have lawfully purchased copies of books then negotiated a license to  
7 reproduce them. Alas, Meta did not even bother to pay the purchase price for the books it illegally  
8 downloaded, let alone obtain a license for their reproduction.

9 45. Meta’s commercial copying of Plaintiff’s work and works owned by the proposed  
10 Class was manifestly unfair use, for several reasons. Many AI companies have described the AI  
11 training process as “teaching” a model human language, much in the way a human learns. While  
12 this self-serving anthropomorphizing of LLMs is misplaced, at a minimum, humans who learn  
13 from books buy them, borrow them from libraries that buy them, or otherwise procure them  
14 lawfully, thus providing at least some measure of compensation to authors and creators. Meta  
15 does not, and it has usurped authors’ content for the purpose of creating a machine built to  
16 generate the very type of content for which authors would usually be paid.

17 46. Meta, in taking authors’ works without compensation, has deprived authors of  
18 books sales and licensing revenues. There is, and has been, an established market for the sale of  
19 books and e-books, yet Meta ignored it and chose to scrape a massive corpus of copyrighted  
20 books from the internet, without even paying for an initial copy.

21 47. Meta has also usurped a licensing market for copyright owners. In the last two  
22 years, a thriving licensing market for copyrighted training data has developed. AI companies have  
23 paid hundreds of millions of dollars to obtain licenses to reproduce high-quality copyrighted  
24 material for LLM training.

25 48. Meta chose to use Plaintiff’s works, and the works owned by the proposed Class,  
26 free of charge, and in doing so has harmed the market for the copyrighted works by depriving  
27 them of book sales and licensing revenue.  
28

1 **III. Meta Pirated Material for Commercial Gain at the Expense of Authors.**

2 49. Perversely, LLMs seriously threaten the livelihood of the same authors whose  
3 works they are non-consensually “trained” on.

4 50. Goldman Sachs’s estimates that generative AI could replace 300 million full-time  
5 jobs in the near future, or one-fourth of labor currently performed in the United States and  
6 Europe.

7 51. Already, writers report losing income from copywriting, journalism, and online  
8 content writing, which are critical sources of income for book authors.

9 52. For example, The Authors Guild, the oldest professional organization in the U.S.  
10 representing writers and authors, and one on the forefront of efforts to shore up “creative markets  
11 against disruptions from generative AI”<sup>25</sup> has published an earnings study that shows a median  
12 writing-related income for full-time authors of just over \$20,000, and that full-time traditional  
13 authors earn only half of that from their books.<sup>26</sup> The rest comes from activities like content  
14 writing—work that is itself starting to dry up as a result of generative AI systems trained on those  
15 writers’ works, without compensation.

16 **IV. Meta Exploited Plaintiff’s Copyrighted Works.**

17 53. Plaintiff and Class Members have suffered identical harms from Meta’s  
18 infringement. The contents of the datasets that Meta used to “train” its LLMs are peculiarly  
19 within its knowledge, such that Plaintiff is unable to discern those contents with perfect accuracy.  
20 But Meta has admitted to using Books3 during the relevant time, and the contents of Books3 is  
21 widely reported. Plaintiff makes specific allegations of infringement below based on what is  
22 known about Meta’s practices and what is known about the contents, uses, and availability of  
23 pirated book repositories that it is suspected Meta used, like Bibliotik.

24  
25  
26 <sup>25</sup> Authors Guild, “Positions and Policy Recommendations”, available at  
<https://authorsguild.org/advocacy/artificial-intelligence/> (last accessed Sept. 26, 2024).

27 <sup>26</sup> Authors Guild, “Top Takeaways from the 2023 Author Income Survey (2023)”, available at  
28 <https://authorsguild.org/news/key-takeaways-from-2023-author-income-survey/#:~:text=Though%20overall%20author%20incomes%20are,coming%20in%20a%20close%20second> (last accessed Sept. 25, 2024).

1 54. Plaintiff Christopher Farnsworth is the author of a number of books, including  
2 *Blood Oath*, *Flashmob: A Novel*, *The Eternal World: A Novel*, and *The President's Vampire*. Each  
3 of these works was and is a part of the Books3 dataset. Pirated copies of these works are  
4 available online through websites like LibGen and Bibliotik. Farnsworth is the author and owner  
5 of the registered copyrights listed under his name in Exhibit A.

6 **CLASS ALLEGATIONS**

7 55. This action is brought by Plaintiff individually and on behalf of the Class, as  
8 defined below, pursuant to Rule 23(a), (b)(3) and 23(b)(2), (c)(4), and (g) of the Federal Rules of  
9 Civil Procedure:

10 All legal and beneficial owners of copyrighted works that: (a) are registered with  
11 the United States Copyright Office; (b) were or are used by Meta in the process of  
12 LLM training, research, or development, including but not limited to the training  
13 and development of its Llama models and (c) have been assigned an International  
Standard Book Number (ISBN). The Class excludes Defendant, its officers and  
directors, members of their immediate families, their co-conspirators, aiders and  
abettors, and the heirs, successors or assigns of any of the foregoing.

14 56. The Class consists of at least thousands of authors and copyright holders and thus  
15 is so numerous that joinder of all members is impractical. The identities of members of the Class  
16 can be readily ascertained from business records maintained by Defendant.

17 57. The claims asserted by Plaintiff are typical of the claims of the Class, all of whose  
18 works were also copied as part of the LLM training, research, and development process.

19 58. The Plaintiff will fairly and adequately protect the interests of the Class and does  
20 not have any interests antagonistic to those of other members of the Class.

21 59. The Plaintiff has retained attorneys who are knowledgeable and experienced in  
22 complex litigation, and have brought multiple copyright class actions against AI companies  
23 asserting that their use of copyrighted material to train LLMs constitutes infringement.

24 60. Plaintiff requests that the Court afford Class members notice and the right to opt-  
25 out of any Class certified in this action.

26 61. This action is appropriate as a class action pursuant to Rule 23(b)(3) of the Federal  
27 Rules of Civil Procedure because common questions of law and fact affecting the Class  
28 predominate over those questions affecting only individual members. The law is uniform. And,

1 the common factual questions giving rise to common answers that move this litigation forward  
2 include:

- 3 a. Whether Meta's reproduction of the Class's copyrighted works constituted  
4 copyright infringement;
- 5 b. Whether Meta's reproduction of the Class's copyrighted works harmed Class  
6 members and whether Class members are entitled to damages, including statutory  
7 damages and the amount of statutory damages; and
- 8 c. Whether Meta's infringement was willful.

9 62. In addition, the class device is the superior mechanism for handling this action,  
10 and a class trial is eminently manageable.

11 63. This action is also appropriate as a class action pursuant to Rule 23(b)(2) of the  
12 Federal Rules of Civil Procedure because Meta's unlicensed exploitation of a large trove of the  
13 Class's books affects all class members in the same way, and any injunctive relief awarded will  
14 affect the Class as a whole.

15 64. Finally, at the very minimum, there are multiple common issues relating to Meta's  
16 uniform conduct, such as (but not limited to) their ingestion, reproduction, and willfulness.

17 **CLAIM FOR RELIEF**  
18 **Copyright Infringement (17 U.S.C. § 501)**  
19 **Against Defendant Meta Platforms, Inc.**

20 65. Plaintiff incorporates by reference the preceding factual allegations.

21 66. Plaintiff and members of the proposed Class have created literary works that are  
22 original and fixed in a tangible medium of expression, and they own the registered copyrights in  
23 the works that Meta reproduced and appropriated for its artificial intelligence projects.

24 67. Plaintiff and members of the proposed Class therefore hold the exclusive rights,  
25 including the rights of reproduction and distribution, to those works under 17 U.S.C. § 106.

26 68. Meta infringed the exclusive rights, under 17 U.S.C. § 106, of Plaintiff and  
27 members of the proposed Class by, among other things, reproducing and/or distributing the works  
28 owned by Plaintiff and the proposed Class in connection with procuring and using datasets for  
artificial intelligence training and development.



1 jury trial for all claims so triable.

2

3 Dated: October 1, 2024

LIEFF CABRASER HEIMANN & BERNSTEIN, LLP

4

5

By: /s/ Daniel Hutchinson

6

Daniel M. Hutchinson

7

Elizabeth J. Cabraser (State Bar No. 083151)

Daniel M. Hutchinson (State Bar No. 239458)

8

Reilly T. Stoler (State Bar No. 310761)

**LIEFF CABRASER HEIMANN**

9

**& BERNSTEIN, LLP**

275 Battery Street, 29th Floor

10

San Francisco, CA 94111-3339

Telephone: (415) 956-1000

11

ecabraser@lchb.com

dhutchinson@lchb.com

12

rstoler@lchb.com

13

Rachel Geman (*pro hac vice forthcoming*)

**LIEFF CABRASER HEIMANN**

14

**& BERNSTEIN, LLP**

250 Hudson Street, 8th Floor

15

New York, New York 10013-1413

Telephone: (212) 355-9500

16

rgeman@lchb.com

17

Scott J. Sholder (*pro hac vice forthcoming*)

CeCe M. Cole (*pro hac vice forthcoming*)

18

**COWAN DEBAETS ABRAHAMS**

**& SHEPPARD LLP**

19

60 Broad Street, 30th Floor

New York, New York 10004

20

Telephone: (212) 974-7474

ssholder@cdas.com

21

ccole@cdas.com

22

*Attorneys for Plaintiff and the Proposed Class*

23

24

25

26

27

28