

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

Elizabeth Brannen (SBN 226234)  
ebrannen@stris.com  
John Stokes (SBN 310847)  
jstokes@stris.com  
Lauren Martin (SBN 294367)  
lmartin@stris.com  
**STRIS & MAHER LLP**  
17785 Center Court Dr N, Ste 600  
Cerritos, CA 90703  
T: (213) 995-6800  
F: (213) 261-0299

Christopher M. Rigali (*pro hac vice*  
forthcoming)  
crigali@stris.com  
Jacqueline Sahlberg (*pro hac vice*  
forthcoming)  
jsahlberg@stris.com  
**STRIS & MAHER LLP**  
1717 K St NW Ste 900  
Washington, DC 20006  
T: (202) 800-5749

*Counsel for Plaintiff*

Kyle Roche (*pro hac vice* forthcoming)  
kroche@fnf.law  
Devin (Velvel) Freedman (*pro hac vice*  
forthcoming)  
vel@fnf.law  
Alex Potter (*pro hac vice* forthcoming)  
apotter@fnf.law  
**FREEDMAN NORMAND FRIEDLAND  
LLP**  
155 E. 44<sup>th</sup> Street, Ste 915  
New York, NY 10017  
T: (646) 494-2900

**UNITED STATES DISTRICT COURT  
NORTHERN DISTRICT OF CALIFORNIA**

COGNELLA, INC.,  
Plaintiff,  
v.  
META PLATFORMS, INC.,  
Defendant.

Civil Case No.:  
**COMPLAINT  
DEMAND FOR JURY TRIAL**

1           1. Plaintiff Cognella, Inc. (“Cognella”) brings this action against Meta Platforms, Inc.  
2 (“Meta” or “Defendant”), and alleges as follows:

3 **I. INTRODUCTION**

4           2. This case concerns Defendant Meta’s exploitation of Plaintiff Cognella’s  
5 copyrighted works to build its “Llama” family of large language models (“LLMs”). Defendant is in  
6 open and active competition with its tech company rivals to win what many have deemed the  
7 “generative AI arms race.” Participants in this race are under immense pressure to build more, better,  
8 and faster, with an eye towards ensuring that *its* generative AI models are widely adopted. Broad  
9 adoption of a company’s generative AI models will yield great profits, the thinking goes; on the  
10 converse, failure to move quickly means missing the boat on this trillion-dollar industry. To remain  
11 competitive in this race, Meta turned to notorious online “shadow libraries” and similar pirated  
12 datasets to obtain vast quantities of copyrighted materials—books, articles, training materials, and  
13 the like—which it desperately wanted to train and optimize its LLM models. Meta could have—and  
14 *should have*—paid copyright owners, including Plaintiff, for licenses to use their copyrighted works  
15 in connection with training its Llama models. Instead, Meta fed from, and then back into, the dark  
16 market for digitally available copyrighted materials. Ultimately, Meta downloaded, reproduced,  
17 distributed, and defaced these copyrighted works, including Plaintiff’s works, to get ahead in the  
18 generative AI arms race and improve its bottom line. Plaintiff brings this action to hold Meta  
19 accountable for its brazen acts of copyright infringement.<sup>1</sup>

20           3. LLMs are a form of “generative artificial intelligence” or “generative AI.” These  
21 models are designed to process and emit natural language. Over the last several years, many of the  
22 world’s biggest and wealthiest technology companies have entered the market to develop, distribute,  
23 and commercialize generative AI models, believing these and other forms of artificial intelligence  
24 have massive potential for growth in revenue and profits. Because artificial intelligence, including  
25  
26

---

27 <sup>1</sup> Unless otherwise indicated, references to Meta’s “Llama” refers to all versions of Llama in any  
28 stage of their development or deployment.

1 generative AI, is viewed by industry leaders as the next foundational layer of the digital economy,  
2 the pressure to build more, better, and faster has led to the aforementioned AI “arms race.”

3 4. Against this backdrop, LLM developers needed data upon which to train their  
4 models. AI researchers quickly discovered that published content—*e.g.*, books, articles, and training  
5 materials—is really the gold standard when it comes to LLM training material. Unlike Internet  
6 content and many other forms of text, published material typically is structured, long-form, and  
7 highly polished. These materials embody the expressive output of their creators, are thoughtfully  
8 planned, and take weeks, months, and years to develop and publish. In simple terms, if you want to  
9 teach a machine how to “speak” or write like a human—capable of telling stories, using analogies,  
10 and making jokes—feed it material that mimics how we express ourselves in our most articulate and  
11 complete forms.

12 5. For its Llama training datasets, Meta needed gold standard training data. To get its  
13 hands on published content, Meta considered and explored the possibility of licensing copyrighted  
14 materials for use in its LLM training. The company’s head of generative AI even discussed spending  
15 up to \$100 million on licensing such materials. Ultimately, however, Meta abandoned this idea in  
16 favor of a faster and cheaper approach: downloading copyrighted materials from Internet-based  
17 “shadow libraries” and other datasets that contained the contents of these shadow libraries. These  
18 shadow libraries, each of which contained tens or hundreds of thousands of unauthorized digital  
19 copies of copyrighted works, had for years been the subject of criminal prosecutions, civil lawsuits,  
20 and widespread warnings within the technology industry. As such, by the time Meta turned to these  
21 sources to obtain training material for its Llama models, it knew or should have known that these  
22 “libraries” were the beneficiaries and propagators of copyright infringement.

23 6. To download copyrighted works from these shadow libraries, Meta often torrented  
24 them. “Torrenting” works by breaking a file into many small pieces and distributing those pieces  
25 across a network of participating computers. A user downloads portions of a file from numerous  
26 other computers that already possess the file and software reassembles those pieces into a completed  
27 whole on the user’s machine.

28

1           7. But torrenting protocols not only facilitate downloading, they also facilitate and  
2 encourage uploading. And through Meta’s torrenting of copyrighted material from illicit shadow  
3 libraries, Meta became a distributor of unauthorized copies of protected works. Stated somewhat  
4 differently, Meta used torrenting protocols that facilitated its reuploading of copyrighted materials  
5 into peer-to-peer file-sharing networks. In Copyright Act terms, Meta distributed copyrighted  
6 materials without consent or authorization. And through this distribution, Meta also contributed to  
7 further acts of infringement, allowing other users on these peer-to-peer networks to unlawfully  
8 reproduce and distribute copyrighted works.

9           8. In addition to ripping pirated copies, reproducing, and distributing without  
10 authorization Plaintiff’s copyrighted materials, Meta also stripped these materials of copyright  
11 management information, enabling, facilitating, and concealing the infringement of these works.  
12 Meta did this to optimize its models’ performance.

13           9. Adding insult to all of this injury, Meta’s use of Plaintiff’s copyrighted materials  
14 facilitated Meta’s creation of AI models capable of generating content that directly competes with  
15 and will compete directly with Plaintiff’s content. “Learning from” the creativity and expression  
16 embodied in thousands upon thousands of copyrighted works, including Plaintiff’s works, Llama  
17 has the capability to flood the market with free and paid content that vies with Plaintiff for consumer  
18 attention.

19           10. The result of Meta’s use of copyrighted materials to train Llama? Meta expects its  
20 Llama models to generate somewhere in the realm of \$460 billion to \$1.4 trillion in revenue. The  
21 Copyright Act doesn’t allow Meta to get a free ride on the backs of Plaintiff and the other creators  
22 whose copyrighted works it exploited to build its trillion-dollar generative AI models. Plaintiff  
23 brings this action to hold Meta accountable for the infringement that enabled its rise in the  
24 generative-AI marketplace, and to enforce the fundamental principle that creative expression cannot  
25 be taken, copied, or exploited without permission or compensation.

26           11. To redress Meta’s repeated, unlawful, and *en masse* infringement of its works,  
27 Plaintiff seeks (1) damages, (2) permanent injunctive relief barring Meta’s ongoing infringement,  
28 and (3) any additional remedies the law provides.

1           12. Plaintiff elects not to bring this case as a class action because the Copyright Act  
2 entitles it to recover individualized statutory damages, determined by a jury, for Meta’s infringement  
3 and related conduct. Plaintiff desires to retain full control of its case and avoid having its rights  
4 diluted by being swept into sprawling class-action settlements structured to resolve claims for  
5 pennies on the dollar. Recent history has shown that certain class actions and proposed settlement(s)  
6 seem to serve the tech conglomerate-infringers, not creators. LLM companies should not be able to  
7 so easily extinguish thousands upon thousands of high-value claims at bargain-basement rates,  
8 eliding what should be the true cost of their massive willful infringement.

9           13. That is not how Plaintiff plans to proceed. Under established Supreme Court  
10 precedent, “the amount of statutory damages is a question for the jury.”<sup>2</sup> The Copyright Act thus  
11 vests authors with the right to have a jury evaluate the willfulness of infringement and assign a  
12 damages amount tailored to the Meta’s conduct.

13           14. In sum, the Copyright Act’s statutory-damages and attorneys’ fee regime empowers  
14 individual authors and publishers to hold infringers accountable without the need for class action  
15 treatment. That is what Plaintiff has chosen to do.

## 16 **II. PARTIES**

### 17 **A. Plaintiff**

18           15. Plaintiff Cognella, Inc., an academic publisher, is a California stock corporation and  
19 the owner or exclusive licensee of thousands of copyrights in textbooks, course readers, custom  
20 course packs, and other learning materials.

21           16. A non-exhaustive list of registered copyrights owned by Plaintiff is included as  
22 Exhibit A (herein, the “Cognella Works”).<sup>3</sup>

23  
24  
25  
26 <sup>2</sup> *Feltner v. Columbia Pictures Television, Inc.*, 523 U.S. 340, 353 (1998).

27 <sup>3</sup> Even where other individuals or entities are listed as copyright claimants on the relevant copyright  
28 registrations, Plaintiff is the owner of all copyrights listed in Exhibit A.

1           **B. Defendant**

2           17. Defendant Meta Platforms, Inc. is a Delaware corporation with its principal place of  
3 business in Menlo Park, California. Meta develops and distributes the Llama series of LLMs,  
4 including Llama 1, Llama 2, Llama 3, and Llama 4, which were trained using datasets sourced in  
5 part from shadow libraries and other datasets containing pirated books.

6           **III. JURISDICTION AND VENUE**

7           18. This action arises under the Copyright Act of 1976, 17 U.S.C. § 101 *et seq.* This  
8 Court has subject-matter jurisdiction under 28 U.S.C. §§ 1331 and 1338(a) because Plaintiff asserts  
9 claims exclusively under federal copyright law.

10           19. This Court has personal jurisdiction over the Defendant. The Defendant has  
11 purposefully availed itself of the privilege of conducting business in this District and the State of  
12 California. The Defendant committed acts of copyright infringement in this District, directed  
13 conduct toward this District, or knowingly caused harm that was suffered in this District. The  
14 Defendant maintains substantial, continuous, and systematic contacts with this District.

15           20. Venue is proper in this District under 28 U.S.C. § 1400(a) because the Defendant or  
16 its agents resides or may be found in this District as a result of the infringing acts alleged herein.  
17 Venue is also proper under 28 U.S.C. § 1391(b)(2) because a substantial part of the events giving  
18 rise to Plaintiff's claims—including the acquisition of pirated copies of Plaintiff's works, the  
19 reproduction and ingestion of those copies into Defendant's training pipelines, the training and fine-  
20 tuning of the relevant LLMs, and the commercialization of the resulting models—occurred in this  
21 District.

22           **IV. FACTUAL ALLEGATIONS**

23           **A. Cognella and its Protected Works**

24           21. Plaintiff Cognella, Inc. is an award-winning academic publisher that produces  
25 instructor-driven, student-centric print and digital learning materials for the higher education  
26 market. For over 30 years, Cognella has developed an expansive catalog and published thousands  
27 of learning materials for instructors and students.

28

1 22. Cognella works with renowned and widely recognized instructors and professionals  
2 across numerous disciplines, many of whom are tenured professors at leading higher education  
3 institutions across the country. In 2025 alone, Cognella had a presence at approximately 1,500  
4 colleges and universities.

5 23. Cognella’s experienced publishing professionals work closely with these subject  
6 matter experts to distill their decades of research and experience into carefully crafted learning  
7 materials with rich pedagogical features and frameworks, transforming their ideas into innovative  
8 educational resources for students, instructors, and the greater educational market. By embracing  
9 authors’ unique perspectives and diverse specialties, Cognella produces timely and valuable titles  
10 that reflect the current needs of the academic market, as well as often underrepresented voices and  
11 topics.

12 24. Cognella offers authors a robust suite of publishing support, including copyediting,  
13 proofreading, layout, design, peer review, developmental consultation, permissions clearance of all  
14 third-party content, developmental assistance of ancillary material, and national marketing and  
15 distribution. This model emphasizes close collaboration with authors, dedication to care and quality,  
16 and the creation of learning materials that reflect the knowledge and skillsets students and  
17 professionals need to be successful in their chosen careers and beyond.

18 25. Cognella works are consistently recognized for their excellence by myriad  
19 organizations and award programs, include the American Journal for Nursing Book of the Year  
20 Awards, the Capstone International Nursing Book Awards, the Graphic Design USA American In-  
21 House Design Awards, the Textbook and Academic Authors Association Awards, the Literary Titan  
22 Book Awards, and the BookFest Awards.

23 26. As a publisher, Cognella has legal or beneficial ownership of the rights to reproduce,  
24 adapt, and distribute thousands of works.

25 27. According to publicly available metadata, the Cognella Works are contained in  
26 pirated online libraries and datasets including Books3 (and thus *The Pile*), Library Genesis or  
27 “LibGen,” Z-Library, and Anna’s Archive. As alleged below, Meta has directly or indirectly  
28 downloaded works contained in *The Pile*, LibGen, Z-Library, and Anna’s Archive, and there is

1 accordingly a reasonable inference that Meta illegally downloaded and reproduced the Cognella  
2 Works.

3 **B. LLMs and the Generative AI “Arms Race”**

4 28. “Generative artificial intelligence” or “generative AI” refers to systems and models  
5 that create outputs—such as text or images—that simulate human expression, often in response to  
6 user prompts.

7 29. Large language models are a form of generative AI, which are designed to process—  
8 or “understand”—and generate natural language. At a high level, LLMs operate by studying and  
9 “learning” the statistical relationship between words and other text (such as punctuation marks);  
10 upon further refinement by their developers, the models then process complex mathematical  
11 sequences to “predict” sequences of text based on the statistical relationships it has learned and been  
12 trained to recognize.

13 30. Development of LLMs requires, among other things, “training” the model on data—  
14 *i.e.*, the model’s inputs. Training “requires identifying a formal measure or ‘objective’ for how well  
15 the model performs, and then repeatedly adjusting the model’s parameters based on that objective  
16 as the model is exposed to training data.”<sup>4</sup> While not all LLM developers use the same terminology,  
17 developers commonly reference a “pre-training” process—“in which a massive amount of  
18 computing power and data is spent to teach the model the broad foundations of language, grammar,  
19 and reasoning”—and a “post-training” or “fine-tuning” process “where the pre-trained model is  
20 further trained on a (relative to pre-training) smaller amount of carefully curated data of specific  
21 tasks.”<sup>5</sup> For purposes of this Complaint, “training” refers to all stages of the LLM training process.

22 31. LLMs typically are trained, at least initially, by feeding them massive amounts of  
23 text data from which they can “learn” statistical relationships between and among text. In the process  
24 of compiling training data and using it during the training process, LLM developers as a matter of  
25

26 <sup>4</sup> U.S. Copyright Office, *Copyright and Artificial Intelligence, Part 3: Generative AI Training Pre-*  
27 *Publication*, at 17 (May 2025), <https://perma.cc/EY5U-EFUY>.

28 <sup>5</sup> *Id.*

1 course typically create multiple copies of “raw” datasets. So, for instance, if an LLM developer  
2 downloads a copy of a digital book from some Internet-based source (e.g., a shadow library), the  
3 developer typically will create a new “cleaned” copy of the book (e.g., stripped of certain  
4 undesirable data), deduplicate that book against other potential copies, compile that book with other  
5 data to create a new compilation (or compilations), and store that book in multiple locations. In  
6 short, with respect to training data, the LLM development pipeline typically involves numerous  
7 reproductions of any given piece of the dataset.

8 32. A model’s inputs—the training datasets—are directly tied to the model’s  
9 performance. As one AI researcher has now famously said, a “model[’s] behavior is not determined  
10 by architecture, hyperparameters, or optimizer choices. It’s determined by your dataset, nothing  
11 else.”<sup>6</sup>

12 33. In determining the corpus of training data for an LLM, a model’s developer typically  
13 considers “the quantity of data, its quality, and the ultimate purpose(s) of the model.”<sup>7</sup>

14 34. In terms of data quantity, the LLMs that have been developed and are being  
15 developed by leading technology companies are trained on enormous datasets, typically on the scale  
16 of terabytes. This is because AI researchers have found that “increasing the quantity of training data  
17 typically increases a model’s ‘performance.’”<sup>8</sup>

18 35. In terms of data quality, “[r]ecent research from major developers suggests that  
19 quality may even be a more important consideration than quantity.”<sup>9</sup> “Garbage in, garbage out,” the  
20 saying goes. AI researchers quickly discovered that published content—books in particular—makes  
21 for some of the best training material for LLMs. Unlike much content on the Internet and many  
22 other forms of text, published material typically is structured, long-form, and highly polished.

---

24 <sup>6</sup> James Betker, *The “it” in AI Models is the Dataset*, (June 10, 2023), <https://perma.cc/ZCH9-H53S>.

25 <sup>7</sup> U.S. Copyright Office, *Copyright and Artificial Intelligence, Part 3: Generative AI Training Pre-*  
26 *Publication*, at 9.

27 <sup>8</sup> *Id.* at 10.

28 <sup>9</sup> *Id.* at 11.

1 Published materials provide formal, extended prose that teaches models narrative structure, complex  
2 syntax, and coherent storytelling.

3 36. Among other sources, academic content is extremely valuable as training material  
4 for LLMs. This is evidenced by the fact that a number of the leading AI companies have been willing  
5 to pay for it from other publishers.<sup>10</sup> For example, Microsoft recently paid Informa, the parent  
6 company of Taylor & Francis, an initial fee of \$10 million to make use of its content “to help  
7 improve relevance and performance of AI systems.”<sup>11</sup> Another leading academic publisher, Wiley,  
8 recently agreed to sell academic content for training AI models by completing a “GenAI content  
9 rights project” with an undisclosed “large tech company” in 2024.<sup>12</sup> Wiley indicated that  
10 negotiations for another generative AI project with a second large tech company were contemplated  
11 for 2025.<sup>13</sup>

12 37. As alleged below, Meta’s employees reached the same conclusion as other AI  
13 researchers—that published materials, and books in particular, were an essential ingredient for  
14 developing a high-performing LLM.

15 38. As discussed in the introduction, technology companies have treated generative AI  
16 as the next foundational layer of the digital economy. Industry leaders publicly describe an “AI arms  
17 race,” in which they have redirected their corporate strategies to seize control of what they believe  
18  
19  
20

---

21 <sup>10</sup> Roger C. Schonfeld, *Tracking the Licensing of Scholarly Content to LLMs*, The Scholarly Kitchen  
22 (Oct. 15, 2024) <https://perma.cc/9DX9-QEWW> (“A number of companies are in the hunt for this  
23 content, including not only OpenAI and Google, but also Apple and more specialized providers.  
24 Investment has been pouring in as a result of the market’s spike in interest in artificial intelligence.”).

25 <sup>11</sup> Christa Dutton, *Two Major Academic Publishers Signed Deals With AI Companies. Some  
26 Professors Are Outraged*, The Chronicle of Higher Education, (July 29, 2024)  
27 <https://perma.cc/PQV5-FWRB>.

28 <sup>12</sup> *Id.*; Press Release, WILEY, *Wiley Increases Quarterly Dividend for the 31<sup>st</sup> Consecutive Year* (June  
27, 2024) <https://perma.cc/JS5D-3WFN>.

<sup>13</sup> Dutton, *Two Major Academic Publishers Signed Deals With AI Companies. Some Professors Are  
Outraged*.

1 will become a new infrastructure layer for commerce, communication, and knowledge work.<sup>14</sup> For  
2 these companies, staying ahead of competitors is “code red.”<sup>15</sup> Among other things, being too slow  
3 out of the blocks could mean ending up in last place; if by the time you publicly released *your* LLM  
4 model, the public writ large and private industry had already adopted *other companies’* models,  
5 there would be a real risk that your model would be left on the shelf.

6 39. Like other participants in this race, Meta risked falling behind early industry leaders  
7 like OpenAI. Indeed, market analysts consistently reported that Meta had, in fact, fallen behind,  
8 with many pointing to Meta’s intense focus on investment in the “Metaverse.”<sup>16</sup>

9 40. Faced with the risks outlined above, Meta needed to move fast to develop and roll  
10 out Llama, its capstone LLM. And to do that, Meta needed to get its hands on high-quality training  
11 data.

### 12 C. Shadow Libraries, *The Pile*, and Torrenting

13 41. For years now, there have been illicit, online marketplaces for digital copies of books.  
14 These “shadow libraries,” as they are commonly known, systematically acquire, store, index, and  
15 disseminate full-fidelity digital copies of copyrighted books. These online repositories maintain  
16 searchable catalogs, metadata, mirrors, bulk-download mechanisms, and—critically—complete  
17 downloadable archives designed to allow third parties to replicate the entire collection.<sup>17</sup> And they

---

18  
19  
20 <sup>14</sup> See Dr. Peter Asaro, *What is an ‘Artificial Intelligence Arms Race’ Anyway?*, 15 I/S: J.L. & Pol’y  
for Info. Soc’y 45 (2019).

21 <sup>15</sup> See Sharon Goldman, *Sam Altman declares ‘Code Red’ as Google’s Gemini surges—three years*  
22 *after ChatGPT cause Google CEO Sundar Pichai to do the same*, FORTUNE (Dec. 2, 2025, 11:43  
AM), <https://perma.cc/J9MS-UQD2>.

23 <sup>16</sup> Mike Proulx, *Meta: So Long, Metaverse; Hello, Superintelligence*, FORRESTER (July 30, 2025),  
24 <https://perma.cc/BZ57-CR2D>; Karen Hao, et al., *Mark Zuckerberg Was Early in AI. Now Meta Is*  
25 *Trying to Catch Up.*, WALL STREET JOURNAL (June 17, 2023, 12:00 AM), [https://perma.cc/G548-](https://perma.cc/G548-B5RD)  
[B5RD](https://perma.cc/G548-B5RD).

26 <sup>17</sup> See Letter from Mary E. Rasenberger, CEO, Authors Guild, & Umair Kazi, Dir. of Pol’y &  
27 Advocacy, Authors Guild, to Daniel Lee, Assistant U.S. Trade Representative for Innovation &  
28 Intellectual Prop., Office of the U.S. Trade Representative (Oct. 7, 2022), [https://perma.cc/XM4R-](https://perma.cc/XM4R-NDN3)  
[NDN3](https://perma.cc/XM4R-NDN3).

1 do all of this without authorization from authors or publishers.<sup>18</sup> In short, shadow libraries facilitate  
2 the reproduction and distribution of unauthorized copies of copyrighted works at industrial scale.

3 42. These libraries are widely known within both piracy communities and the technology  
4 sector as illegal sources of copyrighted books. Many have been the subject of criminal prosecutions,  
5 civil injunctions, domain seizures, and formal designation as “notorious piracy markets” by United  
6 States trade authorities. As centralized shadow libraries increasingly faced enforcement actions,  
7 including the seizure of domains, third parties responded by creating full mirrored copies of those  
8 repositories for decentralized redistribution.<sup>19</sup>

9 43. In 2018, an OpenAI employee downloaded pirated copies of books from Library  
10 Genesis, or “LibGen”—a shadow library repeatedly enjoined by federal courts<sup>20</sup>—and used those  
11 books to create two internal datasets OpenAI called “LibGen1” and “LibGen2,” which OpenAI  
12 publicly referred to as “Books1” and “Books2.”<sup>21</sup> OpenAI used those pirated corpora to train GPT-  
13 3, which it released in June 2020 to widespread commercial acclaim. OpenAI’s use of a books  
14 corpora from a pirate library established the model for how to quickly and cheaply obtain published  
15 materials for use as LLM training data.

16 44. Within weeks of GPT-3’s release, an open research collective, EleutherAI, formed  
17 with the express goal of replicating GPT-3’s capabilities and democratizing access to large-scale  
18 language modeling. To do so, EleutherAI assembled and publicly released *The Pile*, a large (800+

19  
20  
21 <sup>18</sup> See Riddhi Setty, *Rampant ‘Shadow Libraries’ Drive Calls for Anti-Piracy Action*, BLOOMBERG  
22 LAW (Oct. 19, 2022, 9:03 AM), <https://perma.cc/F5VH-3BA6>; Claire Woodcock, *‘Shadow  
23 Libraries’ Are Moving Their Pirated Books to the Dark Web After Fed Crackdowns*, VICE (Nov. 30,  
24 2022, 11:38 AM), <https://perma.cc/K9FA-VLPW>.

25 <sup>19</sup> Woodcock, *‘Shadow Libraries’ Are Moving Their Pirated Books to the Dark Web After Fed  
26 Crackdowns*.

27 <sup>20</sup> See *Cengage Learning, Inc. et al. v. Does 1-50 d/b/a Library Genesis*, ECF No. 36, Case No. 23-  
28 cv-8136 (S.D.N.Y. Sept. 24, 2024).

<sup>21</sup> *In re OpenAI, Inc., Copyright Infringement Litig.*, Case No. 1:2025-md-03143 (S.D.N.Y. 2025),  
ECF No. 846 at 9; Ashley Belanger, *OpenAI desperate to avoid explaining why it deleted pirated  
book datasets*, ARSTECHNICA (Dec. 1, 2025), <https://perma.cc/9M7K-8DA9>.

1 gigabyte) general-purpose training dataset expressly designed to be downloaded, used, and  
2 incorporated into LLMs by academic researchers, startups, and commercial AI developers.<sup>22</sup>

3 45. Because OpenAI did not disclose the precise makeup of its training datasets,  
4 members of EleutherAI constructed a book corpus of their own: “Books3,” consisting of  
5 approximately 196,640 books.<sup>23</sup> Books3 comprises about 12 percent (just over 100 gigabytes) of  
6 *The Pile*.<sup>24</sup>

7 46. EleutherAI and the compiler of Books3, Shawn Presser, have confirmed the genesis  
8 of Books3 is the shadow library Bibliotik. Presser has publicly stated that Books3 represents “all of  
9 bibliotik,”<sup>25</sup> and an EleutherAI paper likewise confirms that Books3 “is a dataset of books derived  
10 from a copy of the contents of the Bibliotik private tracker.”<sup>26</sup>

11 47. Bibliotik is a private, invitation-only torrent tracker that has long functioned as a  
12 centralized source of pirated e-books. The repository, which hosts and distributes hundreds of  
13 thousands of copyrighted books, is only accessible with the use of torrenting protocols, which are  
14 explained below.<sup>27</sup>

15 48. Bibliotik has been widely recognized in piracy communities, academic literature, and  
16 AI-research documentation as a shadow library devoted to copyrighted books. Unsurprisingly, then,  
17 *The Pile*’s datasheet acknowledges that “Books3 is almost entirely comprised of copyrighted  
18  
19  
20

---

21 <sup>22</sup> Leo Gao et al., *The Pile: An 800GB Dataset of Diverse Text for Language Modeling*, arXiv, 1  
22 (2020), <https://perma.cc/NHV6-R8YE>.

23 <sup>23</sup> Stella Biderman et al., *Datasheet for the Pile*, arXiv, 8 (2020), <https://perma.cc/7KL2-LTLF>.

24 <sup>24</sup> *Id.*

25 <sup>25</sup> Shawn Presser, X (Oct. 25, 2020 1:32 AM), <https://perma.cc/7WRD-NHRX>.

26 <sup>26</sup> Biderman et al., *Datasheet for the Pile*, arXiv, 8.

27 <sup>27</sup> Ruheni Mathenge, *The 12 Best Private Torrent Sites Still Working in 2026*, PRIVACYSAVVY (last  
28 accessed March 9, 2026), <https://perma.cc/4V8M-3ALY>.

1 works.”<sup>28</sup> Bibliotik’s illicit nature is no secret—it has been openly discussed for years prior to major  
2 tech companies’ use of datasets derived directly from it.

3 49. In an interview with *The Atlantic*, Presser confirmed that the illuminating purpose  
4 behind Books3 was to ensure broad-based access to the tools necessary to create LLMs:

5 [Presser] created Books3 in the hope that it would allow any developer to create  
6 generative-AI tools. “It would be better if it wasn’t necessary to have something like  
7 Books3,” he said. “But the alternative is that, without Books3, only OpenAI can do  
8 what they’re doing.”<sup>29</sup>

9 50. EletheurAI’s compilation and distribution of *The Pile* and Books3 provided “off-the-  
10 shelf” access to a corpus of infringing works that any AI developer could download and immediately  
11 incorporate into a large-scale training pipeline. As a result, Books3’s presence within *The Pile*  
12 facilitated wide downstream distribution and adoption of a corpus derived from a pirate book  
13 library.<sup>30</sup>

14 51. Bibliotik is just one of many Internet-based shadow libraries and, as alleged below,  
15 Meta relied on more than just *The Pile* in connection with its development of Llama.

16 52. As noted above, OpenAI’s developers trained its LLMs on a books corpus derived  
17 from LibGen, one of the largest and longest-running shadow libraries in the world. LibGen hosts  
18 millions of pirated books, academic texts, and scholarly articles.<sup>31</sup> It operates as a centralized  
19 repository offering direct downloads of full-fidelity e-book files. It also distributes its entire database  
20 through bulk archives and torrent files, enabling third parties to download and locally host complete

21 \_\_\_\_\_  
22 <sup>28</sup> Biderman et al., *Datasheet for the Pile*, arXiv, 8.

23 <sup>29</sup> Alex Reisner, *Revealed: The Authors Whose Pirated Books Are Powering Generative AI*, THE  
24 ATLANTIC (Aug. 19, 2023), <https://perma.cc/P25L-S9FC>.

25 <sup>30</sup> See Stella Biderman, *The Pile: An 800GB Dataset of Diverse Text for Language Modeling*,  
ELEUTHERAI (Dec. 31, 2020), <https://perma.cc/JGT9-LLTK>.

26 <sup>31</sup> Office of the U.S. Trade Representative, REVIEW OF NOTORIOUS MARKETS FOR  
27 COUNTERFEITING AND PIRACY, 27 (2024), <https://perma.cc/N4WT-AHJ9> (“Libgen ... hosts a large  
28 number of digital copies of books, manuals, journals, and other works, many of which are  
unauthorized copies of copyright protected content.”).

1 copies of its collection.<sup>32</sup> LibGen has been repeatedly enjoined by federal courts for copyright  
2 infringement and has been designated a “notorious market” by the United States Trade  
3 Representative.<sup>33</sup> Despite enforcement actions, LibGen has remained accessible through shifting  
4 domains, mirrors, and downloadable archives. Its persistence is a function of deliberate  
5 decentralization designed to evade shutdown.<sup>34</sup>

6 53. Z-Library (also known as “B-ok”) is another well-known shadow library that  
7 emerged as an expanded and user-friendly derivative of LibGen. It incorporated large portions of  
8 LibGen’s catalog while adding additional titles, metadata, and interface features.<sup>35</sup> Z-Library  
9 offered premium features—including faster downloads and higher volume limits—in exchange for  
10 payment, operating in effect as a commercial piracy service.<sup>36</sup> In 2022, Z-Library’s domains were  
11 seized by law-enforcement authorities, and its operators were arrested and later indicted for criminal  
12 copyright infringement.<sup>37</sup> These actions confirmed what had long been publicly known: Z-Library  
13 was an illegal piracy operation. The seizure of Z-Library did not eliminate access to its content.  
14 Instead, third parties responded by creating full mirrors of its collection to ensure continued  
15 distribution.<sup>38</sup>

---

17 <sup>32</sup> See Letter from Mary E. Rasenberger, CEO, Authors Guild, & Umair Kazi, Dir. of Pol’y &  
18 Advocacy, Authors Guild, to Daniel Lee, Assistant U.S. Trade Representative for Innovation &  
19 Intellectual Prop., Office of the U.S. Trade Representative at n.5.

20 <sup>33</sup> See Office of the U.S. Trade Representative, 2019 REVIEW OF NOTORIOUS MARKETS FOR  
21 COUNTERFEITING AND PIRACY, 27, <https://perma.cc/22VN-9VD7>.

22 <sup>34</sup> See Andrew Albanese, *Textbook Publishers Sue Notorious ‘Shadow Library’ Libgen*,  
23 PUBLISHERS WEEKLY (Sep. 14, 2023), <https://perma.cc/3NPY-UJCX>.

24 <sup>35</sup> Jordana Rosenfeld, *Z-Library*, ENCYCLOPEDIA BRITANNICA (last accessed March 9, 2026),  
25 <https://perma.cc/4H26-GDCC>.

26 <sup>36</sup> See Masood Farivar, *Two Russian Nationals Charged With Operating E-Book Piracy Site*, VOA  
27 (Nov. 16, 2022), <https://perma.cc/6MPD-QNKB>.

28 <sup>37</sup> Press Release, U.S. Dep’t of Justice, U.S. Att’y’s Off., E. Dist. of N.Y., *Two Russian Nationals  
Charged with Running Massive E-Book Piracy Website* (Nov. 16, 2022), [https://perma.cc/CR4L-  
JLA3](https://perma.cc/CR4L-JLA3).

<sup>38</sup> See Woodcock, *‘Shadow Libraries’ Are Moving Their Pirated Books to The Dark Web*.

1           54. Another major player in the shadow library ecosystem is Anna’s Archive. Anna’s  
2 Archive is the most comprehensive and active shadow library currently in operation.<sup>39</sup>

3           55. Anna’s Archive began in 2022 as “Pirate Library Mirror,” initially hosting a mirrored  
4 copy of Z-Library. It later rebranded as “Anna’s Archive” and expanded to aggregate and host the  
5 complete collections of LibGen, Z-Library, and other pirated sources.<sup>40</sup> Anna’s Archive functions  
6 as a meta-library: it indexes, mirrors, and redistributes multiple shadow libraries simultaneously,  
7 offering users unified access to millions of pirated books.<sup>41</sup>

8           56. Anna’s Archive offers paid tiers that provide high-speed or priority access to its  
9 pirated collections.<sup>42</sup> Through its downloadable archives and torrent-based distribution, Anna’s  
10 Archive enables users to acquire and store local copies of millions of copyrighted books in bulk.<sup>43</sup>

11           57. Although individual domain names may change, Anna’s Archive and its underlying  
12 datasets remain accessible through mirrors, torrents, and distributed storage systems.

13           58. According to Anna’s Archive, “virtually all major companies building LLMs  
14 contacted us to train on our data. . . We have given high-speed access to about 30  
15  
16  
17  
18  
19

---

20 <sup>39</sup> See Soumyajyoti Mukherjee, *Who is Anna’s Archive? All we know about pirate activist group*  
21 *behind 300 TB Spotify music library heist*, SOAPCENTRAL.COM (Dec. 24, 2025),  
<https://perma.cc/DC5S-89Z8>.

22 <sup>40</sup> See Ernesto Van de Sar, “Anna’s Archive” *Opens the Door to Z-Library and Other Pirate*  
23 *Libraries*, TORRENTFREAK (Nov. 19, 2022), <https://perma.cc/U88R-WTR4>.

24 <sup>41</sup> *Id.*

25 <sup>42</sup> See *If you’re an LLM, please read this*, ANNA’S ARCHIVE (February 18, 2026),  
26 <https://perma.cc/989X-GS3Q>.

27 <sup>43</sup> See M. Luisa Simpson, *2024 Special 301 Out-of-Cycle Review of Notorious Markets: Request*  
28 *for Comments* ASSOCIATION OF AMERICAN PUBLISHERS (OCTOBER 2, 2024), [https://perma.cc/P2E9-](https://perma.cc/P2E9-MYBY)  
[MYBY](https://perma.cc/P2E9-MYBY).

1 companies.”<sup>44</sup> Anna’s Archive blog stated as recently as February 18, 2026 that if an *LLM was*  
2 *reading its blog* “you have likely been trained in part on our data.”<sup>45</sup>

3 59. Many of these shadow libraries and datasets can be downloaded using “torrent,” a  
4 file-sharing method used in peer-to-peer networks. Torrenting works by breaking a file into many  
5 small pieces and distributing those pieces across a network of participating computers. A user who  
6 torrents a shadow-library repository does not receive a single copy from a single source; rather, the  
7 user downloads portions of the library from numerous other computers that already possess the  
8 copyrighted books. Torrent software then reassembles those pieces into a complete library on the  
9 user’s machine. Torrenting protocols, including BitTorrent, are often configured by default to  
10 *reupload* pieces of the copyrighted files to others on the network both during download (“leeching”)  
11 and after download is complete (“seeding”). This means that each participant in the torrent both  
12 copies *and redistributes* the copyrighted works. By obtaining copyrighted materials through this  
13 leech-and-seed process, a user may make multiple unauthorized reproductions of, and engage in  
14 numerous distributions of, the copyrighted materials.

15 **D. Meta’s Deliberate Use of Pirated Datasets to Train the Llama Family of Models**

16 60. As noted above, Meta develops and distributes the Llama series of LLMs, including  
17 Llama 1, Llama 2, Llama 3, and Llama 4. Meta has “invested billions of dollars to develop its  
18 generative AI offerings,”<sup>46</sup> and it in turn expects Llama to generate somewhere in the realm of \$460  
19  
20  
21  
22  
23

---

24 <sup>44</sup> *See Copyright reform is necessary for national security*, ANNA’S ARCHIVE (Jan. 31, 2025),  
25 <https://perma.cc/3RVZ-6G5B>.

26 <sup>45</sup> *See If you’re an LLM, please read this*, ANNA’S ARCHIVE (Feb. 18, 2026), <https://perma.cc/MJ7Z-3ZCL>.

27 <sup>46</sup> *Cambronne Inc., et al. v. Anthropic PBC, et al.*, Case No. 5:2025-cv-10897-PCP (N.D. Cal.), ECF  
28 No. 109 at ¶ 40 (Meta’s Answer).

1 billion to \$1.4 trillion in revenue.<sup>47</sup> None of the money that Meta invested into the development of  
2 its Llama models was used to pay for Plaintiff’s copyrighted works.<sup>48</sup>

3 61. Meta released the first of its Llama models, Llama 1, in February 2023.

4 62. In connection with the development of Llama, Meta viewed book-corpora as among  
5 its most valuable sources of training data. Llama’s design goal was to emit particularly creative and  
6 expressive language, leveraging Meta’s consumer platforms to “connect” with users through text.<sup>49</sup>  
7 To accomplish that, Meta needed to train on large quantities of high-quality books and other  
8 published materials.

9 63. Meta employees repeatedly acknowledged the importance of books as training data.  
10 Specifically, they acknowledged it was “really important for [Meta] to get books data ASAP,” and  
11 the “best resources [Meta] [could] think of are definitely books.”<sup>50</sup>

12 64. Prior to the release of Llama 1, Meta downloaded copyrighted works from LibGen.  
13 As recognized in *Kadrey v. Meta Platforms, Inc.*, Meta first downloaded LibGen in 2022 “to  
14 investigate whether there was value in training Llama on the works it contained. If the answer was  
15 yes, the plan was to then set up licensing agreements for those or similar works.”<sup>51</sup>

16 65. At the time Meta downloaded LibGen, it knew or should have known that LibGen  
17 was a repository of unauthorized copies of copyrighted works.

18 66. Ultimately, Meta decided not to use its pirated LibGen content in connection with  
19 the training of Llama 1. It did, however, rely on another dataset of pirated works—Books3—in  
20 training Llama 1.

21 \_\_\_\_\_  
22 <sup>47</sup> See *Kadrey v. Meta Platforms, Inc.*, 788 F. Supp. 3d 1026, 1040 (N.D. Cal. 2025).

23 <sup>48</sup> *Cambronne*, ECF No. 109 at ¶¶ 2, 9, 42, 89 (admitting in *Cambronne* that “it did not obtain a  
24 license or pay for the use of [copyrighted] works”).

25 <sup>49</sup> Jon Russell, *Mark Zuckerberg Announces New Team at Meta Working on A.I. Products for*  
26 *Instagram, WhatsApp*, CNBC (Feb. 27, 2023), <https://perma.cc/JKE9-N4T7> (“Zuckerberg said that  
the team would build ‘creative and expressive’ tools to be used inside Meta’s products.”).

27 <sup>50</sup> *Kadrey*, 788 F. Supp. 3d at 1040.

28 <sup>51</sup> *Id.*

1           67. In connection with Llama 1’s release, Meta disclosed the composition of the model’s  
2 “pre-training data.”<sup>52</sup> Among the datasets used to train Llama 1 were “two book corpora,” including  
3 “the Books3 section of ThePile.”<sup>53</sup> In their research paper, Meta’s developers described Books3 as  
4 “a publicly available dataset for training large language models,” which they distinguished from the  
5 other book corpus it used, which they described as being in the “pubic domain.”<sup>54</sup> This distinction  
6 was deliberate and notable, as Meta *knew* that Books3 was *not* a corpus of books in the “public  
7 domain.”

8           68. Among the other datasets used by Meta to train Llama 1 were “five CommonCrawl  
9 dumps, ranging from 2017 to 2020,” comprising 67 percent of the overall data mix, and “the publicly  
10 available C4 dataset,” comprising 15 percent of the data mix.<sup>55</sup> “CommonCrawl is a publicly-  
11 available web archive that provides ‘web extracted text’ by removing markup and other non-text  
12 content from the scraped HTML files. This process produces around 20TB of scraped data per  
13 month.”<sup>56</sup> The C4 dataset itself is derived from CommonCrawl.<sup>57</sup> In short, a massive amount (and  
14 the majority) of data used by Meta to train Llama was “scraped” from the Internet.

15           69. As alluded to above, Meta at one time had a “plan . . . to then set up licensing  
16 agreements” to obtain the copyrighted materials it wanted for use in training its AI models.<sup>58</sup> But  
17 this idea was short-lived. “[I]n spring 2023, after failing to acquire licenses and following  
18

---

19 <sup>52</sup> Hugo Touvron et al., *LLaMA: Open and Efficient Foundation Language Models*, arXiv, 2 (2023),  
20 <https://perma.cc/LB97-ZK2C>.

21 <sup>53</sup> *Id.*

22 <sup>54</sup> *Id.*

23 <sup>55</sup> Touvron et al., *LLaMA: Open and Efficient Foundation Language Models*, arXiv, 2.

24 <sup>56</sup> Colin Raffel et al., *Exploring the Limits of Transfer Learning with a Unified Text-to-Text*  
25 *Transformer*, *Journal of Machine Learning Research* 21, 7 (2020), <https://perma.cc/A4VM-V8XR>.

26 <sup>57</sup> *Id.*

27 <sup>58</sup> *Kadrey*, 788 F. Supp. 3d at 1040; *see also Cambronne*, ECF No. 109 at ¶ 88 (Meta admitting that  
28 “it contacted certain publishers” and “internally discussed licensing certain types of data to train its  
Llama models”).

1 escalations to CEO Mark Zuckerberg, Meta decided to just use the works acquired from LibGen as  
2 training data.”<sup>59</sup>

3 70. Done with the idea of paying for licenses to use copyrighted materials, Meta pivoted  
4 to shadow libraries and web crawlers to obtain the training data it wanted. Meta turned to shadow  
5 libraries knowing that these were just that—repositories of unauthorized copies of copyrighted  
6 works.

7 71. “In early 2024, Meta also downloaded Anna’s Archive, a compilation of shadow  
8 libraries including LibGen, Z-Library, and others.”<sup>60</sup>

9 72. Meta likewise downloaded copyrighted materials from Z-Library, Sci-Hub, and  
10 other shadow libraries.<sup>61</sup> Indeed, by April 2024, Meta had downloaded over 25 TB of data from Z-  
11 Library and 10 TB of data from LibGen.<sup>62</sup>

12 73. To obtain the massive quantities of data from these shadow libraries, Meta used  
13 torrenting protocols.<sup>63</sup> In using these protocols to obtain data from these shadow libraries, Meta  
14 reuploaded (and distributed) to peer-to-peer networks at least some of the copyrighted materials it  
15 had downloaded, whether through “leeching” or by “seeding.”<sup>64</sup>

16 74. By uploading and distributing copyrighted materials on peer-to-peer sharing  
17 networks, Meta made those materials available for download to other third parties whose  
18

---

19 <sup>59</sup> *Id.* at 1041.

20 <sup>60</sup> *Kadrey*, 788 F. Supp. 3d at 1041.

21 <sup>61</sup> *Cambronne*, ECF 109 (Meta Answer), ¶ 83.

22 <sup>62</sup> *Kadrey*, ECF No. 472 (Pls.’ Partial Mot. for Summ. J.), at 11.

23 <sup>63</sup> *Cambronne*, ECF No. 109 at ¶ 83 (admitting that “it used BitTorrent to download certain portions  
24 of these publicly available datasets”).

25 <sup>64</sup> *See Kadrey*, 788 F. Supp. 3d at 1041 (“To download these large datasets more quickly and without  
26 unnecessarily slowing down its networks, Meta torrented them. . . . There is no dispute that Meta  
27 torrented LibGen and Anna’s Archive, but the parties dispute whether and to what extent Meta  
28 uploaded (via leeching or seeding) the data it torrented. A Meta engineer involved in the torrenting  
wrote a script to prevent seeding, but apparently not leeching.”).

1 downloading and redistribution of copyrighted materials would constitute further infringement of  
2 the copyrighted materials.

3 75. There was no substantial or commercially significant non-infringing use of the  
4 copyrighted materials that Meta uploaded and distributed. Similarly, there was no substantial or  
5 commercially significant non-infringing use of Meta's uploading and redistribution of these  
6 materials.

7 76. In torrenting from shadow libraries, Meta knew or should have known that it was  
8 participating in a *peer-to-peer network that traded in unauthorized copies of copyrighted material*.  
9 Stated differently, Meta knew or should have known that it was facilitating further copyright  
10 infringement by making available to others on the peer-to-peer network unauthorized copies of  
11 copyrighted materials.

12 77. In connection with the development of Llama, certain Meta employees questioned,  
13 implicitly or explicitly, the legal or ethical implications of relying on copyrighted content, content  
14 from shadow libraries, and torrenting protocols. On information and belief, certain employees  
15 implicitly or explicitly raised concerns with training Llama on "pirated" content. On information  
16 and belief, certain employees implicitly or explicitly raised concerns with obtaining copyrighted  
17 materials from shadow libraries. And on information and belief, certain employees implicitly or  
18 explicitly raised concerns with torrenting from shadow libraries.

19 78. Notwithstanding what it knew about these peer-to-peer networks and concerns raised  
20 by employees, Meta torrented large quantities of data from these shadow libraries, including but not  
21 limited to terabytes of data from each of Anna's Archive and LibGen.

22 79. Meta continued to download from these shadow libraries even after it had been sued  
23 for copyright infringement based on its use of copyrighted works to train Llama.

24 80. Meta also relied on web crawlers to obtain massive amounts of Internet data for use  
25 in training Llama. On information and belief, the web crawler data that Meta obtained in connection  
26 with its Llama training contained copyrighted materials.

27 81. Upon downloading this material, Meta manipulated the downloaded data. Among  
28 other things, Meta processed the data, and in doing so, removed and altered certain text and

1 information, including titles, author information, information about the copyright owner, and  
2 copyright notices. When it removed this information, Meta did so without the authority of the  
3 copyright owners.

4 82. As Meta has admitted elsewhere, Meta “created a script to remove repetitive text  
5 from the training data for its Llama models.”<sup>65</sup> As admitted elsewhere, Meta’s “program was  
6 designed to remove lines that contain the word ‘copyright.’”<sup>66</sup>

7 83. As alleged above, the Cognella Works appear in various shadow libraries and/or  
8 pirated datasets, including Books3, LibGen, Anna’s Archive, and Z-Library. As such, there is a  
9 reasonable inference that Meta downloaded, manipulated, reproduced, and distributed the Cognella  
10 Works.

11 84. As admitted elsewhere, Meta used at least Books3 and data it obtained from LibGen  
12 and Z-Library to train the Llama family of models.<sup>67</sup> On information and belief, it also used data it  
13 obtained from Anna’s Archive and other shadow libraries to train the Llama family of models. After  
14 Llama 1, Meta stopped disclosing detailed information about the datasets that its models were  
15 trained on. Meta did this, in part, to conceal the fact that it had relied on and was relying on pirated  
16 materials to train Llama.

17 85. On information and belief, Meta, like Anthropic, downloaded, maintained, and  
18 retained pirated copies of works that were not specifically tied to training its LLM models.<sup>68</sup>

19 86. Meta has incorporated Llama into several of its consumer-facing products, including  
20 Facebook, Instagram, and WhatsApp. And although Meta has offered free downloads of various  
21 versions of Llama, it considers or has considered Llama and the technology underpinning it to be

22 \_\_\_\_\_  
23 <sup>65</sup> *Kadrey*, ECF 485, ¶ 88 (Meta’s Answer to Third Am. Compl.).

24 <sup>66</sup> *Kadrey*, ECF 485, ¶ 89.

25 <sup>67</sup> *Cambronne*, ECF 109, ¶ 2 (Meta’s Answer, admitting that “it obtained publicly available data  
26 from sources such as LibGen and Z-Library, portions of which were utilized to research, develop,  
and train its Llama models”).

27 <sup>68</sup> *Bartz v. Anthropic PBC*, 787 F. Supp. 3d 1007, 1014 (N.D. Cal. 2025) (Anthropic’s amassing of  
28 “a central library”).

1 an important part of its commercial success. In other words, Meta hopes and expects that its  
2 development of Llama will contribute to its current and future financial success.

3 **E. Meta’s Conduct Harms the Market for Plaintiff’s Copyrighted Materials**

4 87. Meta’s conduct has a detrimental effect on the potential market for and value of  
5 Plaintiff’s works, including by, among other things, developing products that create and are capable  
6 of creating content which serves as a direct substitute for Plaintiff’s works, developing products that  
7 create content and are capable of creating content which serves as indirect substitutes for Plaintiff’s  
8 works, and undermining Plaintiff’s ability to participate in and profit from the market for licensing  
9 its works for the purpose of training LLMs.

10 88. *First*, Meta’s decision to download and use unauthorized copies of Plaintiff’s works  
11 from shadow libraries deprived Plaintiff of revenue in the form of licensing fees that it would have  
12 otherwise earned. As alleged herein, there is an existing market for licensing copyrighted materials  
13 such as Plaintiff’s, including for use in the development of LLMs. Meta bypassed that market, and  
14 in doing so, deprived Plaintiff of licensing revenue it would have earned. Meta now allows users in  
15 certain markets to opt out of having their data used to train its AI models, but it deprived Plaintiff  
16 and many others of that choice. Meta’s misconduct undermined Plaintiff and many others,  
17 eliminating the bargaining power they should have had, and otherwise would have had, with respect  
18 to licensing terms for the use Meta made of their works.

19 89. *Second*, Llama, trained on Plaintiff’s copyrighted materials (and the protected works  
20 of others), is capable of generating outputs that compete directly with, and risk serving as  
21 replacements for, Plaintiff’s works. As the U.S. Copyright Office has warned, “the speed and scale  
22 at which AI systems generate content pose a serious risk of diluting markets for works of the same  
23 kind as in their training data.”<sup>69</sup>

24 90. *Third*, even if Llama models are restricted from outputting extended portions of  
25 verbatim text from copyrighted works, they are nevertheless capable of producing nearly  
26

---

27 <sup>69</sup> U.S. Copyright Office, *Copyright and Artificial Intelligence, Part 3: Generative AI Training Pre-*  
28 *Publication*, at 65.

1 indistinguishable “versions” of copyrighted works such that a consumer would use the AI-generated  
2 version of the material rather than pay for a copy of the actual copyrighted work.

3 **V. CLAIMS FOR RELIEF**

4 **COUNT I**

5 **Direct Copyright Infringement (17 U.S.C. § 501)**

6 91. Plaintiff incorporates the allegations above.

7 92. Plaintiff is the legal or beneficial owner of the copyrighted works listed in Exhibit A  
8 (referred to herein as the Cognella Works).

9 93. Defendant, without authorization from Plaintiff, copied, downloaded, reproduced,  
10 ingested, parsed, embedded, and used pirated copies of Plaintiff’s works in the development,  
11 training, fine-tuning, and deployment of its Llama family of models. These acts violated Plaintiff’s  
12 exclusive rights under 17 U.S.C. § 106.

13 94. Defendant’s infringement occurred repeatedly throughout the lifecycle of its AI-  
14 model development. As alleged above, Defendant:

- 15 • acquired Plaintiff’s works from shadow-library repositories and datasets containing
- 16 pirated works from shadow libraries;
- 17 • distributed Plaintiff’s works through the use of torrenting software, programs, or
- 18 protocols;
- 19 • reproduced additional copies during ingestion, preprocessing, storage, deduplication,
- 20 formatting, and/or tokenization; and
- 21 • while training its models, made even more copies of the text—because every training
- 22 pass (each epoch and each step of gradient descent) automatically requires creating
- 23 and working with fresh versions of that text.

24 95. Defendant’s reproductions and distributions of Plaintiff’s copyrighted works were  
25 made without permission, license, or consent and violated Plaintiff’s exclusive rights under the  
26 Copyright Act.

27 96. Defendant’s infringement was **willful**. As alleged above, Defendant knowingly  
28 trained its models on and/or optimized its product with datasets saturated with pirated books,

1 including Plaintiff's works; relied on shadow-library corpora it knew to be illegal; ignored internal  
2 and external warnings; attempted to conceal the composition of its training datasets; and continued  
3 copying after public reports, lawsuits, law-enforcement seizures, cease-and-desist notices, and  
4 industry-wide alerts made the illegality unmistakable.

5 97. Upon information and belief, Defendant has made and will continue to make  
6 substantial profits and gains to which it is not in law or in equity entitled.

7 98. Plaintiff has been injured by Defendant's willful acts of copyright infringement.  
8 Plaintiff is entitled to statutory damages, actual damages, restitution of profits, and/or other remedies  
9 in law or equity.

10 99. Plaintiff is entitled to recover attorneys' fees and costs under 17 U.S.C. § 505.

11 **COUNT II**

12 **Contributory Copyright Infringement (17 U.S.C. § 501)**

13 100. Plaintiff incorporates the allegations above.

14 101. Defendant used torrenting software, programs, or protocols to download datasets  
15 containing pirated copies of works, including the Cognella Works.

16 102. In connection with its torrenting of datasets that contained copyrighted works,  
17 Defendant uploaded and distributed, either through "seeding" and/or "leeching," copyrighted  
18 materials, including the Cognella Works, thereby making those works available to third parties for  
19 downloading on peer-to-peer networks.

20 103. Defendant knowingly participated in peer-to-peer sharing networks that it knew  
21 trafficked in pirated copies of copyrighted materials. In other words, Defendant knew that others on  
22 these networks were infringing copyrighted materials through reproduction and/or distribution.  
23 There was no substantial or commercially significant non-infringing use of the copyrighted  
24 materials that Meta uploaded and distributed. Nor was there substantial or commercially significant  
25 non-infringing use of Defendant's uploading and distribution of Plaintiff's copyrighted works. By  
26 participating in these networks, and by further uploading and distributing Plaintiff's copyrighted  
27 works, Defendant materially contributed to and induced further infringement of Plaintiff's works.

28

1 104. By knowingly inducing and materially contributing to others' infringement of  
2 Plaintiff's works, Defendant is liable for contributory copyright infringement.

3 105. As a direct and proximate cause of Defendant's conduct, Plaintiff was injured and is  
4 entitled to statutory damages, actual damages, restitution of profits, and/or other remedies in law or  
5 equity.

6 106. Plaintiff is entitled to recover attorneys' fees and costs under 17 U.S.C. § 505.

7 **COUNT III**

8 **Removal of Copyright Management Information (17 U.S.C. § 1202(b)(1))**

9 107. Plaintiff incorporates the allegations above.

10 108. Plaintiff's works contain information that constitutes "copyright management  
11 information" as that term is defined in 17 U.S.C. § 1202(c). This includes but is not limited to author  
12 information, information about the copyright owner, and copyright notices.

13 109. Upon downloading copyrighted materials, including Plaintiff's works, Defendant  
14 processed the data, and in doing so, removed and altered certain text and information, including  
15 copyright management information found in and on Plaintiff's works. When it removed this  
16 information, Defendant did so without the authority of the copyright owners.

17 110. Defendant's removal of the copyright management information was intentional—it  
18 did so to, among other things, create high-quality LLM training data and, through the creation and  
19 use of high-quality training data, ultimately create high-quality LLM models.

20 111. Defendant removed copyright management information from Plaintiff's works  
21 knowing or having reasonable grounds to believe that it was enabling, facilitating, and concealing  
22 acts of copyright infringement. As to concealment, Defendant knew or had reasonable grounds to  
23 believe that by stripping copyrighted works of copyright management information it would be  
24 harder for others to discover the true sources—*e.g.*, copyrighted works—of Defendant's training  
25 data.

26 112. Plaintiff was harmed by Defendant's removal of copyright management information  
27 from its works and is entitled to statutory damages, actual damages, restitution of profits, and other  
28

1 remedies provided by law. Plaintiff is entitled to recover attorneys' fees and costs under 17 U.S.C.  
2 § 1203(b)(5).

3 **PRAYER FOR RELIEF**

4 WHEREFORE, Plaintiff requests that the Court enter judgment on its behalf by ordering:

- 5 a. Judgment in favor of Plaintiff against the Defendant;
- 6 b. A declaration that the Defendant has infringed Plaintiff's exclusive copyrights  
7 under the Copyright Act;
- 8 c. A declaration that such infringement is willful;
- 9 d. A declaration that Defendant violated 17 U.S.C. § 1202(b) through its removal  
10 of copyright management information;
- 11 e. A permanent injunction enjoining the Defendant and all those acting in concert  
12 with it from engaging in the infringing conduct alleged herein;
- 13 f. That the Defendant be directed to account to Plaintiff for all gains, profits, and  
14 advantages derived from their unlawful acts;
- 15 g. An award of statutory damages under the Copyright Act;
- 16 h. An award of statutory or actual damages under 17 U.S.C. § 1203(c);
- 17 i. An award of restitution, disgorgement, costs, expenses, and attorneys' fees as  
18 permitted by law (including those allowable under 17 U.S.C. § 505 and/or 17  
19 U.S.C. § 1203(b)(4)–(5)).
- 20 j. Pre- and post-judgment interest on the damages awarded to Plaintiff; and
- 21 k. Further relief for Plaintiff as the Court may deem just and proper.

22 **JURY TRIAL DEMANDED**

23 Under Federal Rule of Civil Procedure 38(b), Plaintiff demands a trial by jury.  
24  
25  
26  
27  
28

1 Dated: May 4, 2026

Respectfully submitted,

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

/s/ Elizabeth Brannen

Elizabeth Brannen (SBN 226234)  
John Stokes (SBN 310847)  
Lauren Martin (SBN 294367)  
**STRIS & MAHER LLP**  
17785 Center Court Dr N, Ste 600  
Cerritos, CA 90703  
T: (213) 995-6800  
F: (213) 261-0299  
ebrannen@stris.com  
jstokes@stris.com  
lmartin@stris.com

Christopher M. Rigali (*pro hac vice  
forthcoming*)  
Jacqueline Sahlberg (*pro hac vice forthcoming*)  
1717 K St NW Ste 900  
Washington, DC 20006  
Phone: (202) 800-5749  
crigali@stris.com  
jsahlberg@stris.com

Kyle Roche (*pro hac vice forthcoming*)  
Devin (Velvel) Freedman (*pro hac vice  
forthcoming*)  
Alex Potter (*pro hac vice forthcoming*)  
**FREEDMAN NORMAND FRIEDLAND  
LLP**  
155 E. 44<sup>th</sup> Street, Ste 915  
New York NY 10017  
T: (646) 494-2900  
vel@fnf.law  
kroche@fnf.law  
apotter@fnf.law

*Counsel for Plaintiff*