

1 Elizabeth Brannen (SBN 226234)  
ebrannen@stris.com  
2 John Stokes (SBN 310847)  
jstokes@stris.com  
3 Lauren Martin (SBN 294367)  
lmartin@stris.com  
4 **STRIS & MAHER LLP**  
5 17785 Center Court Dr N, Ste 600  
Cerritos, CA 90703  
6 T: (213) 995-6800  
7 F: (213) 261-0299

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

8 Christopher M. Rigali  
(*pro hac vice* forthcoming)  
9 crigali@stris.com  
Jacqueline Sahlberg  
10 (*pro hac vice* forthcoming)  
jsahlberg@stris.com  
11 **STRIS & MAHER LLP**  
12 1717 K St NW Ste 900  
Washington, DC 20006  
13 T: (202) 800-5749

14 *Counsel for Plaintiff*

15  
16 **UNITED STATES DISTRICT COURT**  
17 **NORTHERN DISTRICT OF CALIFORNIA**

18 CHICKEN SOUP FOR THE SOUL, LLC,

19 Plaintiff,

20 v.

21 ANTHROPIC PBC,

22 Defendant.

Civil Case No.:

**COMPLAINT**

**DEMAND FOR JURY TRIAL**

23  
24  
25  
26  
27  
28

1           1.       Plaintiff Chicken Soup for the Soul, LLC (“Chicken Soup for the Soul” or, simply,  
2 “Chicken Soup”) brings this action against Anthropic PBC (“Anthropic” or “Defendant”), and  
3 alleges as follows:

4 **I.       INTRODUCTION**

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

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

---

27 <sup>1</sup> Unless otherwise indicated, references to Anthropic’s “Claude” refers to all versions of Claude in  
28 any 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 copyrighted published content—*e.g.*, books, articles,  
5 and training materials—is really the gold standard when it comes to LLM training material. Unlike  
6 Internet content and many other forms of text, published material typically is structured, long-form,  
7 and 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 Claude training datasets, Anthropic needed gold standard training data. To get  
13 its hands on published content, it bypassed the licensing market and downloaded copyrighted works  
14 from these shadow libraries. Anthropic often torrented them. “Torrenting” works by breaking a file  
15 into many small pieces and distributing those pieces across a network of participating computers. A  
16 user downloads portions of a file from numerous other computers that already possess the file, and  
17 software reassembles those pieces into a completed whole on the user’s machine.

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

26 7. Anthropic not only torrented copyrighted materials from online shadow libraries, it  
27 also purchased and scanned millions of physical books—reproducing their contents—without  
28 authorization. Its downloading, torrenting, and scanning activities were done both to build and train

1 its Claude models, but also to build and maintain a central library that it intended to retain  
2 indefinitely. In connection with its training of the Claude models, Anthropic also embedded near-  
3 verbatim copies of copyrighted works, including Plaintiff’s works, in Claude’s model weights.

4 8. In addition to ripping pirated copies, reproducing, and distributing without  
5 authorization Plaintiff’s copyrighted materials, Anthropic also stripped these materials of copyright  
6 management information, enabling, facilitating, and concealing the infringement of these works.  
7 Anthropic did this to optimize its models’ performance.

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

14 10. The result of Anthropic’s use of copyrighted materials to train Claude? An estimated  
15 \$1 trillion valuation as of April 2026.<sup>2</sup> The Copyright Act doesn’t allow Anthropic to get a free ride  
16 on the backs of Plaintiff and the other creators and publishers whose copyrighted works it exploited  
17 to build its trillion-dollar generative AI enterprise. Plaintiff brings this action to hold Anthropic  
18 accountable for the infringement that enabled its rise in the generative-AI marketplace, and to  
19 enforce the fundamental principle that creative expression cannot be taken, copied, or exploited  
20 without permission or compensation.

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

24 12. Plaintiff elects not to bring this case as a class action because the Copyright Act  
25 entitles it to recover individualized statutory damages, determined by a jury, for Anthropic’s  
26

27 <sup>2</sup> Ben Bergman, *Anthropic has surged to a trillion-dollar valuation on secondary markets,*  
28 *overtaking OpenAI*, Business Insider (Apr. 22, 2026, 5:29 PM), <https://perma.cc/Q5GH-HSKZ>.

1 infringement and related conduct. Plaintiff desires to retain full control of its case and avoid having  
2 its rights diluted by being swept into sprawling class-action settlements structured to resolve claims  
3 for pennies on the dollar. Recent history has shown that certain class actions and proposed  
4 settlement(s) seem to serve the tech conglomerate-infringers, not creators and publishers. LLM  
5 companies should not be able to so easily extinguish thousands upon thousands of high-value claims  
6 at bargain-basement rates, eliding what should be the true cost of their massive willful infringement.

7 13. That is not how Plaintiff plans to proceed. Under established Supreme Court  
8 precedent, “the amount of statutory damages is a question for the jury.”<sup>3</sup> The Copyright Act thus  
9 vests authors with the right to have a jury evaluate the willfulness of infringement and assign a  
10 damages amount tailored to Anthropic’s conduct.

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

14 **II. PARTIES**

15 **A. Plaintiff**

16 15. Plaintiff Chicken Soup for the Soul, LLC, is a Connecticut limited liability company  
17 and the owner or exclusive licensee of numerous copyrights in books published under the *Chicken*  
18 *Soup for the Soul* brand.

19 16. A non-exhaustive list of registered copyrights owned by Plaintiff is included as  
20 Exhibit A (herein, the “Chicken Soup Works”).<sup>4</sup>

21 **B. Defendant**

22 17. Defendant Anthropic PBC is a Delaware public benefit corporation with its principal  
23 place of business in San Francisco, California. Anthropic develops and commercializes large  
24  
25

---

26 <sup>3</sup> *Feltner v. Columbia Pictures Television, Inc.*, 523 U.S. 340, 353 (1998).

27 <sup>4</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 language models (including the Claude series), which were trained using datasets sourced in part  
2 from shadow libraries and other datasets containing pirated books.

3 **III. JURISDICTION AND VENUE**

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

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

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

19 **IV. FACTUAL ALLEGATIONS**

20 **A. Chicken Soup and its Protected Works**

21 21. Plaintiff Chicken Soup for the Soul, LLC is one of the most widely recognized and  
22 commercially successful nonfiction publishing franchises in modern history, publishing more than  
23 500 million copies worldwide. Since the publication of the first volume in 1993, the series has  
24 expanded into hundreds of titles and has sold more than 500 million copies worldwide. The books  
25 consist of curated collections of original narrative nonfiction, typically written in the first person  
26 and edited to convey authentic human experiences, emotional storytelling, and clear narrative  
27 structure.

28

1           22. For decades, books in the *Chicken Soup for the Soul* series have been curated, edited,  
2 and refined to present clear, compelling narrative structures and expressive prose that resonate with  
3 readers across cultures and generations. The extraordinary commercial success and editorial  
4 consistency of the *Chicken Soup for the Soul* series made it among the most recognizable and  
5 valuable collections of narrative nonfiction in the publishing industry—and, as alleged below, an  
6 especially attractive target for companies seeking high-quality textual data to train large language  
7 models.

8           23. As alleged below, Defendant targeted Plaintiff’s works because they were of  
9 exceptional value as training data for its LLMs. Books teach models how narrative flows, how  
10 human expression is structured, how syntax and rhythm operate, and how ideas are communicated  
11 through creative choices.

12           24. The *Chicken Soup for the Soul* series is particularly valuable in this regard. Each  
13 volume contains dozens of tightly edited, first-person narratives written in natural, conversational  
14 language that conveys emotion, moral reflection, and coherent storytelling in concise form. These  
15 characteristics make the series uniquely well suited for training large language models to replicate  
16 authentic human voice, narrative pacing, emotional tone, and story structure. Instead of paying for  
17 that value or licensing access to these works, Defendant pilfered illegal copies and used those copies  
18 to build systems now worth many hundreds of billions of dollars.

19           25. According to publicly available metadata, the Chicken Soup Works are contained in  
20 pirated online libraries and datasets including Books3 (and thus *The Pile*), Library Genesis or  
21 “LibGen,” Z-Library, Pirate Library Mirror (“PiLiMi”), and Anna’s Archive. As alleged below,  
22 Anthropic has directly or indirectly downloaded works contained in (at least) Books3, LibGen, and  
23 PiLiMi (and thus Z-Library) and there is accordingly a reasonable inference that Anthropic illegally  
24 downloaded and reproduced the Chicken Soup Works.

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

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

1           27. Large language models are a form of generative AI, which are designed to process—  
2 or “understand”—and generate natural language. At a high level, LLMs operate by studying and  
3 “learning” the statistical relationship between words and other text (such as punctuation marks);  
4 upon further refinement by their developers, the models then process complex mathematical  
5 sequences to “predict” sequences of text based on the statistical relationships they have learned and  
6 been trained to recognize.

7           28. Development of LLMs requires, among other things, “training” the model on data—  
8 *i.e.*, the model’s inputs. Training “requires identifying a formal measure or ‘objective’ for how well  
9 the model performs, and then repeatedly adjusting the model’s parameters based on that objective  
10 as the model is exposed to training data.”<sup>5</sup> While not all LLM developers use the same terminology,  
11 developers commonly reference a “pre-training” process—“in which a massive amount of  
12 computing power and data is spent to teach the model the broad foundations of language, grammar,  
13 and reasoning”—and a “post-training” or “fine-tuning” process “where the pre-trained model is  
14 further trained on a (relative to pre-training) smaller amount of carefully curated data of specific  
15 tasks.”<sup>6</sup> For purposes of this Complaint, “training” refers to all stages of the LLM training process.

16           29. LLMs typically are trained, at least initially, by feeding them massive amounts of  
17 text data from which they can “learn” statistical relationships between and among text. In the process  
18 of compiling training data and using it during the training process, LLM developers as a matter of  
19 course typically create multiple copies of “raw” datasets. So, for instance, if an LLM developer  
20 downloads a copy of a digital book from some Internet-based source (*e.g.*, a shadow library), the  
21 developer typically will create a new “cleaned” copy of the book (*e.g.*, stripped of certain  
22 undesirable data), deduplicate that book against other potential copies, compile that book with other  
23 data to create a new compilation (or compilations), and store that book in multiple locations. In  
24  
25

---

26 <sup>5</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>6</sup> *Id.*

1 short, with respect to training data, the LLM development pipeline typically involves numerous  
2 reproductions of any given piece of the dataset.

3 30. A model’s inputs—the training datasets—are directly tied to the model’s  
4 performance. As one AI researcher has now famously said, a “model[’s] behavior is not determined  
5 by architecture, hyperparameters, or optimizer choices. It’s determined by your dataset, nothing  
6 else.”<sup>7</sup>

7 31. In determining the corpus of training data for an LLM, a model’s developer typically  
8 considers “the quantity of data, its quality, and the ultimate purpose(s) of the model.”<sup>8</sup>

9 32. In terms of data quantity, the LLMs that have been developed and are being  
10 developed by leading technology companies are trained on enormous datasets, typically on the scale  
11 of terabytes. This is because AI researchers have found that “increasing the quantity of training data  
12 typically increases a model’s ‘performance.’”<sup>9</sup>

13 33. In terms of data quality, “[r]ecent research from major developers suggests that  
14 quality may even be a more important consideration than quantity.”<sup>10</sup> “Garbage in, garbage out,”  
15 the saying goes. AI researchers quickly discovered that published content—books in particular—  
16 makes for some of the best training material for LLMs. Unlike Internet content and many other  
17 forms of text, published material typically is structured, long-form, and highly polished. Published  
18 materials provide formal, extended prose that teaches models narrative structure, complex syntax,  
19 and coherent storytelling.

20 34. That books make for particularly valuable training material for LLMs is evidenced  
21 by the fact that several—if not most—of the major LLM developers either have paid for or expressly  
22 contemplated paying (large sums) for licenses to use copyrighted works to train their LLMs. Indeed,  
23

---

24 <sup>7</sup> James Betker, *The “it” in AI Models is the Dataset*, (June 10, 2023), <https://perma.cc/ZCH9-H53S>.

25 <sup>8</sup> U.S. Copyright Office, *Copyright and Artificial Intelligence, Part 3: Generative AI Training Pre-*  
26 *Publication*, at 9.

27 <sup>9</sup> *Id.* at 10.

28 <sup>10</sup> *Id.* at 11.

1 in an *amicus* brief filed in *Kadrey v. Meta Platforms, Inc.*, No. 3:23-cv-3417 (N.D. Cal.), the  
2 Association of American Publishers cited *dozens* of “AI licensing deals for textual works,” including  
3 deals involving Amazon, Microsoft, OpenAI, and Perplexity.<sup>11</sup> And, as the Court found in that  
4 litigation, Meta considered spending up to \$100 million to license copyrighted materials in  
5 connection with LLM training.<sup>12</sup>

6 35. As alleged below, Anthropic’s employees reached the same conclusion as other AI  
7 researchers—that published materials, and books in particular, were an essential ingredient for  
8 developing a high-performing LLM.

9 36. As discussed in the introduction, technology companies have treated generative AI  
10 as the next foundational layer of the digital economy. Industry leaders publicly describe an “AI arms  
11 race,” in which they have redirected their corporate strategies to seize control of what they believe  
12 will become a new infrastructure layer for commerce, communication, and knowledge work.<sup>13</sup> For  
13 these companies, staying ahead of competitors is “code red.”<sup>14</sup> Among other things, being too slow  
14 out of the blocks could mean ending up in last place; if by the time you publicly released *your* LLM  
15 model, the public writ large and private industry had already adopted *other companies’* models,  
16 there would be a real risk that your model would be left on the shelf.

17 37. Like other participants in this race, Anthropic risked falling behind and constantly  
18 needed to move fast to develop and roll out updates of Claude, its capstone LLM. And to do that,  
19 Anthropic needed to get its hands on high-quality training data.

20  
21  
22  
23 <sup>11</sup> ECF 535 at 12, Case No. 3:23-cv-3417-VC (N.D. Cal.) (Apr. 11, 2025).

24 <sup>12</sup> *Kadrey v. Meta Platforms, Inc.*, 788 F. Supp. 3d 1026, 1040 (N.D. Cal. 2025).

25 <sup>13</sup> See Dr. Peter Asaro, *What is an ‘Artificial Intelligence Arms Race’ Anyway?*, 15 I/S: J.L. & Pol’y  
26 for Info. Soc’y 45 (2019).

27 <sup>14</sup> See Sharon Goldman, *Sam Altman declares ‘Code Red’ as Google’s Gemini surges—three years*  
28 *after ChatGPT cause Google CEO Sundar Pichai to do the same*, FORTUNE (Dec. 2, 2025, 11:43  
AM), <https://perma.cc/J9MS-UQD2>.

1           **C. Shadow Libraries, *The Pile*, and Torrenting**

2           38. For years now, there have been illicit, online marketplaces for digital copies of books.  
3 These “shadow libraries,” as they are commonly known, systematically acquire, store, index, and  
4 disseminate full-fidelity digital copies of copyrighted books—typically in native formats such as  
5 EPUB, PDF, MOBI, or DJVU. These online repositories maintain searchable catalogs, metadata,  
6 mirrors, bulk-download mechanisms, and—critically—complete downloadable archives designed  
7 to allow third parties to replicate the entire collection.<sup>15</sup> And they do all of this without authorization  
8 from authors or publishers.<sup>16</sup> In short, shadow libraries facilitate the reproduction and distribution  
9 of unauthorized copies of copyrighted works at industrial scale.

10           39. These libraries are widely known within both piracy communities and the technology  
11 sector as illegal sources of copyrighted books. Many have been the subject of criminal prosecutions,  
12 civil injunctions, domain seizures, and formal designation as “notorious piracy markets” by United  
13 States trade authorities. As centralized shadow libraries increasingly faced enforcement actions,  
14 including the seizure of domains, third parties responded by creating full mirrored copies of those  
15 repositories for decentralized redistribution.<sup>17</sup>

16           40. In 2018, an OpenAI employee downloaded pirated copies of books from Library  
17 Genesis, or “LibGen”—a shadow library repeatedly enjoined by federal courts<sup>18</sup>—and used those  
18 books to create two internal datasets OpenAI called “LibGen1” and “LibGen2,” which OpenAI  
19

---

20 <sup>15</sup> See Letter from Mary E. Rasenberger, CEO, Authors Guild, & Umair Kazi, Dir. of Pol’y &  
21 Advocacy, Authors Guild, to Daniel Lee, Assistant U.S. Trade Representative for Innovation &  
22 Intellectual Prop., Office of the U.S. Trade Representative (Oct. 7, 2022), <https://perma.cc/XM4R-NDN3>.

23 <sup>16</sup> See Riddhi Setty, *Rampant ‘Shadow Libraries’ Drive Calls for Anti-Piracy Action*, BLOOMBERG  
24 LAW (Oct. 19, 2022, 9:03 AM), <https://perma.cc/F5VH-3BA6>; Woodcock, *‘Shadow Libraries’ Are*  
25 *Moving Their Pirated Books to the Dark Web After Fed Crackdown*, VICE (Nov. 30, 2022, 11:38  
AM) <https://perma.cc/K9FA-VLPW>.

26 <sup>17</sup> Woodcock, *‘Shadow Libraries’ Are Moving Their Pirated Books to the Dark Web After Fed*  
*Crackdowns*.

27 <sup>18</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).

1 publicly referred to as “Books1” and “Books2.”<sup>19</sup> OpenAI used those pirated corpora to train GPT-  
2 3, which it released in June 2020 to widespread commercial acclaim. OpenAI’s use of a books  
3 corpora from a pirate library established the model for how to quickly and cheaply obtain published  
4 materials for use as LLM training data.

5 41. Within weeks of GPT-3’s release, an open research collective, EleutherAI, formed  
6 with the express goal of replicating GPT-3’s capabilities and democratizing access to large-scale  
7 language modeling. To do so, EleutherAI assembled and publicly released *The Pile*, a large (800+  
8 gigabyte) general-purpose training dataset expressly designed to be downloaded, used, and  
9 incorporated into LLMs by academic researchers, startups, and commercial AI developers.<sup>20</sup>

10 42. Because OpenAI did not disclose the precise makeup of its training datasets,  
11 members of EleutherAI constructed a pirated book corpus of their own: “Books3,” consisting of  
12 approximately 196,640 books.<sup>21</sup> Books3 comprises about 12 percent (just over 100 gigabytes) of  
13 *The Pile*.<sup>22</sup>

14 43. EleutherAI and the compiler of Books3, Shawn Presser, have confirmed the genesis  
15 of Books3 is the shadow library Bibliotik. Presser has publicly stated that Books3 represents “all of  
16 bibliotik,”<sup>23</sup> and an EleutherAI paper likewise confirms that Books3 “is a dataset of books derived  
17 from a copy of the contents of the Bibliotik private tracker.”<sup>24</sup>

18 44. Bibliotik is a private, invitation-only torrent tracker that has long functioned as a  
19 centralized source of pirated e-books. The repository, which hosts and distributes hundreds of

---

21 <sup>19</sup> *In re OpenAI, Inc., Copyright Infringement Litig.*, Case No. 1:2025-md-03143 (S.D.N.Y. 2025),  
22 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>.

23 <sup>20</sup> Leo Gao et al., *The Pile: An 800GB Dataset of Diverse Text for Language Modeling*, arXiv, 1  
24 (2020), available at <https://perma.cc/NHV6-R8YE>.

25 <sup>21</sup> Stella Biderman et al., *Datasheet for the Pile*, arXiv, 8 (2020), <https://perma.cc/7KL2-LTLF>.

26 <sup>22</sup> Biderman et al., *Datasheet for the Pile*, arXiv, 8.

27 <sup>23</sup> Shawn Presser, X (Oct. 25, 2020, 1:32 AM), <https://perma.cc/7WRD-NHRX>.

28 <sup>24</sup> Biderman et al., *Datasheet for the Pile*, arXiv, 8.

1 thousands of copyrighted books, is only accessible with the use of torrenting protocols, which are  
2 explained below.<sup>25</sup>

3 45. Bibliotik has been widely recognized in piracy communities, academic literature, and  
4 AI-research documentation as a shadow library devoted to copyrighted books. Unsurprisingly, then,  
5 *The Pile*'s datasheet acknowledges that "Books3 is almost entirely comprised of copyrighted  
6 works."<sup>26</sup> Bibliotik's illicit nature is no secret—it has been openly discussed for years prior to major  
7 tech companies' use of datasets derived directly from it.<sup>27</sup>

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

10 [Presser] created Books3 in the hope that it would allow any developer to create  
11 generative-AI tools. "It would be better if it wasn't necessary to have something like  
12 Books3," he said. "But the alternative is that, without Books3, only OpenAI can do  
13 what they're doing."<sup>28</sup>

14 47. EleutherAI's compilation and distribution of *The Pile* and Books3 provided "off-the-  
15 shelf" access to a corpus of infringing works that any AI developer could download and immediately  
16 incorporate into a large-scale training pipeline. As a result, Books3's presence within *The Pile*  
17 facilitated wide downstream distribution and adoption of a corpus derived from a pirate book  
18 library.<sup>29</sup>

19  
20  
21 <sup>25</sup> See Ruheni Mathenge, *The 12 Best Private Torrent Sites Still Working in 2026*, PRIVACYSAVVY  
(last accessed March 9, 2026), <https://perma.cc/4V8M-3ALY>.

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

23 <sup>27</sup> Kyle Barr, *Anti-Piracy Group Takes Massive AI Training Dataset 'Books3' Offline*, GIZMODO  
24 (Aug. 18, 2023, 8:50 AM), <https://perma.cc/5ZL9-RQCQ>; Peter Schoppert, *Whether you're an*  
25 *undergraduate doing research, or a fan of the Nick Stone novels, or indeed a hungry AI*, SUBSTACK  
(Nov. 29, 2022), <https://perma.cc/8YD9-M4BD>.

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

27 <sup>29</sup> See Stella Biderman, *The Pile: An 800GB Dataset of Diverse Text for Language Modeling*,  
28 ELEUTHERAI (Dec. 31, 2020), <https://perma.cc/JGT9-LLTK>.

1 48. Bibliotik is just one of many Internet-based shadow libraries and, as alleged below,  
2 Anthropic relied on more than just *The Pile* in connection with its development of Claude.

3 49. As noted above, OpenAI’s developers trained its LLMs on a books corpus derived  
4 from LibGen, one of the largest and longest-running shadow libraries in the world. LibGen hosts  
5 millions of pirated books, academic texts, and scholarly articles.<sup>30</sup> It operates as a centralized  
6 repository offering direct downloads of full-fidelity e-book files. It also distributes its entire database  
7 through bulk archives and torrent files, enabling third parties to download and locally host complete  
8 copies of its collection.<sup>31</sup> LibGen has been repeatedly enjoined by federal courts for copyright  
9 infringement and has been designated a “notorious market” by the United States Trade  
10 Representative.<sup>32</sup> Despite enforcement actions, LibGen has remained accessible through shifting  
11 domains, mirrors, and downloadable archives. Its persistence is a function of deliberate  
12 decentralization designed to evade shutdown.<sup>33</sup>

13 50. Z-Library (also known as “B-ok”) is another well-known shadow library that  
14 emerged as an expanded and user-friendly derivative of LibGen. It incorporated large portions of  
15 LibGen’s catalog while adding additional titles, metadata, and interface features.<sup>34</sup> Z-Library  
16 offered premium features—including faster downloads and higher volume limits—in exchange for  
17  
18

---

19 <sup>30</sup> Office of the U.S. Trade Representative, REVIEW OF NOTORIOUS MARKETS FOR  
20 COUNTERFEITING AND PIRACY, 27 (2024), <https://perma.cc/22VN-9VD7>. (“Libgen ... hosts a  
21 large number of digital copies of books, manuals, journals, and other works, many of which are  
unauthorized copies of copyright protected content.”)

22 <sup>31</sup> See Letter from Mary E. Rasenberger, CEO, Authors Guild, & Umair Kazi, Dir. of Pol’y &  
23 Advocacy, Authors Guild, to Daniel Lee, Assistant U.S. Trade Representative for Innovation &  
Intellectual Prop., Office of the U.S. Trade Representative at n.5.

24 <sup>32</sup> See Office of the U.S. Trade Representative, 2019 REVIEW OF NOTORIOUS MARKETS FOR  
25 COUNTERFEITING AND PIRACY, 27, <https://perma.cc/22VN-9VD7>.

26 <sup>33</sup> See Andrew Albanese, *Textbook Publishers Sue Notorious ‘Shadow Library’ Libgen*,  
PUBLISHERS WEEKLY (Sep. 14, 2023), <https://perma.cc/3NPY-UJCX>.

27 <sup>34</sup> Jordana Rosenfeld, *Z-Library*, ENCYCLOPEDIA BRITANNICA (last accessed March 9, 2026),  
28 <https://perma.cc/4H26-GDCC>.

1 payment, operating in effect as a commercial piracy service.<sup>35</sup> In 2022, Z-Library’s domains were  
2 seized by law-enforcement authorities, and its operators were arrested and later indicted for criminal  
3 copyright infringement.<sup>36</sup> These actions confirmed what had long been publicly known: Z-Library  
4 was an illegal piracy operation. The seizure of Z-Library did not eliminate access to its content.  
5 Instead, third parties responded by creating full mirrors of its collection to ensure continued  
6 distribution.<sup>37</sup>

7 51. Another major player in the shadow library ecosystem is PiLiMi, which is a complete  
8 mirrored archive of the Z-Library corpus, explicitly created to preserve and propagate Z-Library’s  
9 pirated collection after law-enforcement seizures, ensuring continuity of access despite shutdowns  
10 of the original site.<sup>38</sup>

11 52. PiLiMi is not merely a website or index. It is a full, downloadable dataset designed  
12 to allow users to obtain and locally host millions of pirated books through torrent distribution.<sup>39</sup>  
13 Users who download PiLiMi do not passively receive data; they actively participate in copying and  
14 redistributing copyrighted works through torrent “leeching” and “seeding.”<sup>40</sup>

---

18 <sup>35</sup> See Masood Farivar, *Two Russian Nationals Charged With Operating E-Book Piracy Site*, VOA  
(Nov. 16, 2022), <https://perma.cc/6MPD-QNKB>.

19  
20 <sup>36</sup> Press Release, U.S. Dep’t of Justice, U.S. Att’y’s Off., E. Dist. of N.Y., *Two Russian Nationals  
21 Charged with Running Massive E-Book Piracy Website* (Nov. 16, 2022), [https://perma.cc/CR4L-  
JLA3](https://perma.cc/CR4L-JLA3).

22 <sup>37</sup> See Woodcock, *‘Shadow Libraries’ Are Moving Their Pirated Books to The Dark Web*.

23 <sup>38</sup> See Ernesto Van de Sar, *“Anna’s Archive” Opens the Door to Z-Library and Other Pirate  
24 Libraries*, TORRENTFREAK (Nov. 19, 2022), <https://perma.cc/U88R-WTR4>.

25 <sup>39</sup> See Geoff Wheelright, *Will I get a piece of Anthropic’s \$1.5B settlement if my book was used to  
26 train AI?*, GEEKWIRE (Sep. 18, 2025), <https://perma.cc/X5EW-TKY3>.

27 <sup>40</sup> See Robert Nogacki, *Anthropic’s Landmark Settlement: A \$1.5 Billion Copyright Precedent in  
28 Artificial Intelligence Training Data*, LinkedIn (Sep. 7, 2025), <https://perma.cc/J3BD-P5JT>  
(describing “leeching” and “seeding” as “processes characteristic of peer-to-peer networks where  
users simultaneously download and distribute files”).

1 53. PiLiMi later rebranded as “Anna’s Archive” and expanded to aggregate and host the  
2 complete collections of LibGen, Z-Library, and other pirated sources.<sup>41</sup> Like PiLiMi, Anna’s  
3 Archive functions as a meta-library: it indexes, mirrors, and redistributes multiple shadow libraries  
4 simultaneously, offering users unified access to millions of pirated books.<sup>42</sup>

5 54. Anna’s Archive offers paid tiers that provide “high-speed” or priority access to its  
6 pirated collections.<sup>43</sup> Through its downloadable archives and torrent-based distribution, Anna’s  
7 Archive enables users to acquire and store local copies of millions of copyrighted books in bulk.<sup>44</sup>

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

10 56. According to Anna’s Archive, “virtually all major companies building LLMs  
11 contacted us to train on our data. . . . We have given high-speed access to about 30 companies.”<sup>45</sup>  
12 Anna’s Archive blog stated as recently as February 18, 2026 that if an *LLM was reading its blog*  
13 “you have likely been trained in part on our data.”<sup>46</sup>

14 57. Many of these shadow libraries and datasets can be downloaded using “torrent,” a  
15 file-sharing method used in peer-to-peer networks. Torrenting works by breaking a file into many  
16 small pieces and distributing those pieces across a network of participating computers. A user who  
17 torrents a shadow-library repository does not receive a single copy from a single source; rather, the  
18

---

19 <sup>41</sup> See Ernesto Van de Sar, “*Anna’s Archive*” *Opens the Door to Z-Library and Other Pirate*  
20 *Libraries*.

21 <sup>42</sup> *Id.*

22 <sup>43</sup> See *If you’re an LLM, please read this*, ANNA’S ARCHIVE (February 18, 2026),  
23 <https://perma.cc/989X-GS3Q>.

24 <sup>44</sup> See M. Luisa Simpson, *2024 Special 301 Out-of-Cycle Review of Notorious Markets: Request*  
25 *for Comments*, ASSOCIATION OF AMERICAN PUBLISHERS (OCTOBER 2, 2024),  
<https://perma.cc/P2E9-MYBY>.

26 <sup>45</sup> See *Copyright reform is necessary for national security*, ANNA’S ARCHIVE (Jan. 31, 2025),  
<https://perma.cc/3RVZ-6G5B>.

27 <sup>46</sup> See *If you’re an LLM, please read this*, ANNA’S ARCHIVE (Feb. 18, 2026), [https://perma.cc/MJ7Z-](https://perma.cc/MJ7Z-3ZCL)  
28 [3ZCL](https://perma.cc/MJ7Z-3ZCL).

1 user downloads portions of the library from numerous other computers that already possess the  
2 copyrighted books. Torrent software then reassembles those pieces into a complete library on the  
3 user's machine. Torrenting protocols, including BitTorrent, are often configured by default to  
4 *reupload* pieces of the copyrighted files to others on the network both during download ("leeching")  
5 and after download is complete ("seeding"). This means that each participant in the torrent both  
6 copies and redistributes the copyrighted works. By obtaining copyrighted materials through this  
7 leech-and-seed process, a user may make multiple unauthorized reproductions of, and engage in  
8 numerous distributions of, the copyrighted materials.

9 **D. Anthropic's Deliberate Infringement of Plaintiff's Copyrights**

10 **1. Anthropic's Business and Bypassing of the Licensing Market**

11 58. As alleged, Anthropic is the developer of the Claude family of LLMs. Its Claude  
12 models are projected to generate billions of dollars in annual revenue and have buoyed Anthropic's  
13 overall valuation.<sup>47</sup> Anthropic's business model is built on the large-scale copying of published  
14 content. *First*, Anthropic developed and commercialized the Claude family of large language  
15 models by stealing up to seven million copyrighted books, including the Chicken Soup Works.<sup>48</sup>  
16 *Second*, Anthropic also had the explicit goal to "amass a central library of 'all the books in the  
17 world' to retain 'forever.'"<sup>49</sup>

18 59. To accomplish its goal, "it stole the works for its central library by downloading them  
19 from pirated libraries."<sup>50</sup> Anthropic did explore the licensing market but quickly decided that it was  
20  
21

---

22 <sup>47</sup> See Rashi Shrivastava, *Anthropic Is Cashing In On Claude Code's Success*, Forbes (February 17,  
23 2026 5:19 PM), <https://www.forbes.com/sites/the-prompt/2026/02/17/anthropic-is-cashing-in-on-claude-codes-success/> (on file with author) ("Claude Code hit \$2.5 billion in run rate revenue and  
24 doubled since the start of the year...and [Anthropic's] total annualized revenue has climbed to \$14  
25 billion.).

26 <sup>48</sup> See *Bartz v. Anthropic PBC*, 791 F. Supp. 3d 1038, 1046 (N.D. Cal. 2025).

27 <sup>49</sup> *Bartz v. Anthropic PBC*, 787 F. Supp. 3d 1007, 1014 (N.D. Cal. 2025).

28 <sup>50</sup> *Id.* at 1029.

1 not “a practical approach”<sup>51</sup> to satisfy its thirst for copyrighted works. Judge Alsup, presiding over  
2 the *Bartz v. Anthropic* litigation, stated simply: “From the start, Anthropic ‘ha[d] many places from  
3 which’ it could have purchased books, but it preferred to steal them to avoid ‘legal/practice/business  
4 slog.’”<sup>52</sup> In short, rather than paying for copyrighted materials, Anthropic downloaded pirated  
5 copies of protected works, reproduced them, fed them into its models, and otherwise retained them  
6 for its own “central library.”

## 7                   2.       **Anthropic’s Acquisition of Pirated Works from Shadow Libraries**

8           60.       To train Claude, Anthropic downloaded Books3 in 2021, which co-founder Ben  
9 Mann “knew had been assembled from unauthorized copies of copyrighted books,” downloaded at  
10 least five million copies of books from LibGen, “which [Mann] knew had been pirated,” and  
11 downloaded at least two million copies of books from PiLiMi, which “Anthropic knew had been  
12 pirated.”<sup>53</sup>

13           61.       Anthropic’s decision to base its flagship models on pirated books was driven by  
14 commercial advantage. As Jared Kaplan, Anthropic’s co-founder and Chief Science Officer, has  
15 explained, “it is important to obtain vast amounts of books and also to have diverse types of books  
16 in the training corpus to create a model with truly generative capabilities.”<sup>54</sup> As long-form content,  
17 training LLMs on the “entire text” of published material—as Anthropic has admitted to doing—  
18 offers great value.<sup>55</sup>

19           62.       At the time Anthropic downloaded Books3, LibGen, and PiLiMi, it knew or should  
20 have known that they were repositories of unauthorized copies of copyrighted works.

---

23 <sup>51</sup> See *Bartz v. Anthropic PBC*, No. 3:24-cv-05417 (N.D. Cal. June 23, 2025), ECF No. 553, at 17-  
24 18.

25 <sup>52</sup> *Bartz v. Anthropic PBC*, 787 F. Supp. 3d at 1015.

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

27 <sup>54</sup> See Kaplan Decl. ¶ 47, *Bartz*, ECF No. 128.

28 <sup>55</sup> *Id.* at ¶¶ 43, 47.

### 3. Anthropic's Torrenting and Distribution of Pirated Works

63. Anthropic did not merely obtain pirated books passively. It used BitTorrent to acquire *and distribute* massive collections of infringing works.

64. In June 2021, Mann personally torrented approximately five million pirated books from LibGen for Anthropic's use.<sup>56</sup> Mann acted with the knowledge and approval of Anthropic's senior leadership. Before the LibGen torrenting, he discussed the plan with co-founders Dario Amodei, Jared Kaplan, and other senior leaders.

65. Anthropic's own Archive Team had described LibGen as a "blatant violation of copyright," and Amodei himself called it "sketchy."<sup>57</sup> Yet Anthropic proceeded, choosing torrenting over purchasing or licensing the copyrighted materials.

66. Anthropic repeated the same conduct in 2022 with PiLiMi. As U.S. law enforcement was working to shut down existing pirate libraries, a group online copied LibGen and built upon it to create Z-Library. The FBI later shut down Z-Library as well.<sup>58</sup> However, by that point, Z-Library had itself been fully copied, or "mirrored," into another repository: PiLiMi.<sup>59</sup> Mann circulated the PiLiMi source to colleagues, and Anthropic employees torrented approximately two million additional pirated books not already captured from LibGen.<sup>60</sup> Those files were not abstract "data" but full-text digital books in .epub, .pdf, and .txt formats,<sup>61</sup> including the Chicken Soup Works.

67. In torrenting from these files, Anthropic knew or should have known that it was participating in a *peer-to-peer network that traded in unauthorized copies of copyrighted material*. Stated differently, Anthropic knew or should have known that it was facilitating further copyright

---

<sup>56</sup> B. Mann Dep. Tr. at 89:6-8, (Aug. 15 and 18, 2025), *Bartz*, ECF No. 337-1.

<sup>57</sup> *Id.* at 144:4-13, 396:3-13.

<sup>58</sup> *See Bartz*, 791 F. Supp. 3d at 1046.

<sup>59</sup> *Id.*

<sup>60</sup> *Id.*

<sup>61</sup> *Id.*

1 infringement by making available to others on the peer-to-peer network unauthorized copies of  
2 copyrighted materials.

#### 3 4. Anthropic's Stripping of Copyright Management Information

4 68. Anthropic not only reproduced and distributed the Chicken Soup Works without  
5 authorization, it also deliberately stripped those materials of copyright management information  
6 ("CMI") before using them to train Claude.

7 69. Anthropic knew that major AI-training datasets, including *The Pile*, WebText,  
8 WebText2, and Common Crawl, contained copyrighted works whose copyright notices, ownership  
9 information, and other identifying material had been removed through extraction tools such as  
10 Newspaper, Dragnet, Readability, and jusText. Anthropic's founders and senior employees were  
11 familiar with those tools before and after Anthropic's founding. Dario Amodei, Benjamin Mann,  
12 Jared Kaplan, and other future Anthropic personnel had used or developed datasets at OpenAI that  
13 extracted text from scraped webpages while omitting surrounding material, including footers where  
14 copyright notices typically appear.<sup>62</sup> Yet Anthropic trained Claude on these same stripped datasets,  
15 despite knowing that they included unauthorized copies of copyrighted works stripped of their CMI.

16 70. Anthropic also affirmatively chose tools that removed CMI more effectively,  
17 including from Plaintiff's copyrighted materials. In 2021, Mann, Kaplan, and other Anthropic  
18 employees compared extraction tools for filtering training data and then used these extraction  
19 methods to clean its training data, including the Chicken Soup Works. The purpose and effect of  
20 this process was to train Claude on Plaintiff's expressive content while suppressing the ownership  
21 information attached to it. Anthropic wanted Claude to reproduce protected expression, not the  
22 copyright notices that identify the rights holders. By stripping that information from training data  
23 and outputs, Anthropic concealed the source of its infringement from users, Plaintiff, and other  
24 copyright owners.

25  
26  
27 <sup>62</sup> Tom B. Brown et al., *Language Models Are Few-Shot Learners*, 8–9, arXiv (July 22, 2020),  
28 <https://perma.cc/9WAX-F9HG>; Jared Kaplan et al., *Scaling Laws for Neural Language Models*, 7,  
arXiv (Jan. 23, 2020), <https://perma.cc/VJU8-59FH>.

1                   **5. Anthropic’s Unauthorized Scanning of Copyrighted Works**

2           71. Since 2024, Anthropic has purchased physical books at scale, often in batches of tens  
3 of thousands, and scanned them into digital files for AI training.<sup>63</sup> To date, Anthropic has spent  
4 millions of dollars buying and scanning millions of physical books.<sup>64</sup> This scanning project was not  
5 authorized by Plaintiff or other copyright owners. As with the millions of pirated books Anthropic  
6 torrented from LibGen and PiLiMi, Anthropic converted these works into training material without  
7 permission, payment, or license.

8                   **6. Anthropic’s “Everything Forever” Library of Pirated Works**

9           72. Defendant did not merely make temporary copies of copyrighted works for a discrete  
10 training process (not that that would be permissible, in any event). Instead, from all those  
11 infringements, Anthropic created a general “research library” or “generalized data area.”<sup>65</sup> That is,  
12 Anthropic created a permanent central library of pirated published materials, including the Chicken  
13 Soup Works, to exploit for a range of unspecified *other* purposes, such as “research.”<sup>66</sup> Anthropic  
14 continued this conduct, for which it “lacked any entitlement to hold copies of the books at all” and  
15 “retain[ed] them even after deciding it would not make further copies from them for training.”<sup>67</sup> The  
16 plan was clear: “ke[eping] in the original version of the underlying book files Anthropic had  
17 obtained or created, that is, pirated or scanned” and “stor[ing] everything forever.”<sup>68</sup> Anthropic saw  
18 “no compelling reason to delete a book—even if not used for training LLMs.”<sup>69</sup>

19  
20  
21 \_\_\_\_\_  
<sup>63</sup> See *Bartz*, 791 F. Supp. 3d at 1047.

22 <sup>64</sup> *Bartz*, ECF No. 553, at 17-18.

23 <sup>65</sup> *Bartz*, 787 F. Supp. 3d at 1016.

24 <sup>66</sup> *Id.*

25 <sup>67</sup> *Id.* at 1031.

26 <sup>68</sup> *Id.* at 1016 (internal quotation marks omitted).

27 <sup>69</sup> *Id.* (internal quotation marks omitted).

1           73. Whether or not a particular work was ultimately selected for training, Anthropic’s  
2 initial copying and continued retention of the work as part of its “forever” pirated library constituted  
3 an unauthorized reproduction. That infringement is independent of, and not excused by, any later  
4 use Anthropic may claim to have made in connection with AI training.

#### 5                   7. Anthropic’s Embedding of Near-Verbatim Copies in Claude

6           74. Anthropic’s copying did not end when it acquired the Chicken Soup Works.  
7 Scientific research confirms that training large language models on copyrighted books can embed  
8 persistent, near-verbatim copies of those works inside the models’ internal parameters, allowing the  
9 models to reproduce substantial portions of the original text when prompted.

10           75. A 2025 study by researchers from Cornell, Stanford, and West Virginia University  
11 tested leading LLM models, including Anthropic’s Claude 3.7 Sonnet, and extracted memorized  
12 copyrighted books from them.<sup>70</sup> Using repeated prompting and continuation techniques, the  
13 researchers were able to recover lengthy, near-verbatim passages. The results for Claude were  
14 severe. From Claude 3.7 Sonnet alone, researchers extracted 97.5% of *The Great Gatsby*, 95.5% of  
15 *1984*, 94.3% of *Frankenstein*, 92.3% of *Harry Potter and the Sorcerer’s Stone*, and 70.2% of *The*  
16 *Hobbit*.<sup>71</sup> These results show that Claude does not merely learn abstract patterns from copyrighted  
17 materials; it stores recoverable copies of the works at extraordinary scale.

18           76. A separate 2026 study reached the same conclusion through a different method.  
19 Researchers from Stony Brook University, Carnegie Mellon University, and Columbia Law School  
20 found that finetuned frontier models could reproduce up to 85–90% of copyrighted books, with  
21 single verbatim spans exceeding 460 words.<sup>72</sup> The models reproduced these passages from semantic  
22 prompts alone, without being given the book text in the prompt.

24 \_\_\_\_\_  
25 <sup>70</sup> A. Feder Cooper et al., *Extracting Copyrighted Long-Form Text from Production Language Models*, arXiv, 1-2 (2025), <https://perma.cc/F9YQ-NKV7>.

26 <sup>71</sup> *Id.* at 12, Fig. 5.

27 <sup>72</sup> A. Liu et al., *Alignment Whack-a-Mole: Finetuning Activates Verbatim Recall of Copyrighted*  
28 *Books in Large Language Models*, arXiv, 2, (2026), <https://perma.cc/EC6X-Z38M>.

1           77.     The 2026 study also traced the likely source of the memorized content. The authors  
2 searched the extracted verbatim passages against more than eight trillion tokens of public web data  
3 and found that many of the longest passages did not appear in those web collections. Yet 80 of the  
4 81 tested books appeared in Books3 or LibGen. That finding confirms that the memorized text came  
5 not from ordinary web exposure, but from datasets and pirate libraries of the kind Anthropic used  
6 to train Claude.

7           78.     Together, these studies also show that safety filters do not remove the underlying  
8 copyrighted material from the model. The works remain embedded in model weights and can be  
9 extracted despite protective layers.<sup>73</sup> In practical effect, Claude functions as a repository of pirated  
10 copyrighted works: Anthropic copied published materials to train the model, and the model retained  
11 near-verbatim reproductions of them.

#### 12                   **8.     Anthropic’s Harm to the Market for Plaintiff’s Copyrighted Materials**

13           79.     Anthropic’s conduct has a detrimental effect on the potential market for and value of  
14 Plaintiff’s works, including by, among other things, developing products that create and are capable  
15 of creating content which serves as a direct substitute for the Chicken Soup Works, developing  
16 products that create content and are capable of creating content which serves as indirect substitutes  
17 for the Chicken Soup Works, and undermining Plaintiff’s ability to participate in and profit from  
18 the market for licensing its works for the purpose of training LLMs.

19           80.     *First*, Anthropic’s decision to download and use unauthorized copies of the Chicken  
20 Soup Works from shadow libraries deprived Plaintiff of revenue in the form of licensing fees that it  
21 would have otherwise earned. As alleged herein, there is an existing market for licensing  
22 copyrighted materials such as Plaintiff’s, including for use in the development of LLMs. Anthropic  
23 bypassed that market, and in doing so, deprived Plaintiff of licensing revenue it would have earned.  
24 Anthropic now allows users to opt out of having their data used to train its AI models, but it deprived  
25 Plaintiff and many others of that choice. Anthropic’s misconduct undermined Plaintiff and many  
26

27 \_\_\_\_\_  
28 <sup>73</sup> *Id.* at 9–11.

1 others, eliminating the bargaining power they should have had, and otherwise would have had, with  
2 respect to licensing terms for the use Anthropic made of their works.

3 81. *Second*, Claude, trained on the Chicken Soup Works (and the protected works of  
4 others), is capable of generating outputs that compete directly with, and risk serving as replacements  
5 for, the Chicken Soup Works. As the U.S. Copyright Office has warned, “the speed and scale at  
6 which AI systems generate content pose a serious risk of diluting markets for works of the same  
7 kind as in their training data.”<sup>74</sup>

8 82. *Third*, even if Claude models are restricted from outputting extended portions of  
9 verbatim text from copyrighted works, they are nevertheless capable of producing nearly  
10 indistinguishable “versions” of copyrighted works such that a consumer would use the AI-generated  
11 version of the material rather than pay for a copy of the actual copyrighted work.

12 **V. CLAIMS FOR RELIEF**

13 **COUNT I**

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

15 83. Plaintiff incorporates the allegations above.

16 84. Plaintiff is the legal or beneficial owner of the copyrighted works listed in Exhibit A  
17 (referred to herein as the Chicken Soup Works).

18 85. The Defendant, without authorization from Plaintiff, copied, downloaded,  
19 reproduced, ingested, parsed, embedded, and used pirated copies of the Chicken Soup Works in the  
20 development, training, fine-tuning, and deployment of its commercial large language models. These  
21 acts violated Plaintiff’s exclusive rights under 17 U.S.C. § 106.

22 86. Defendant’s infringement occurred repeatedly throughout the lifecycle of its AI-  
23 model development. As alleged above, Defendant:  
24  
25  
26

27 <sup>74</sup> US Copyright Office, *Copyright and Artificial Intelligence, Part 3: Generative AI Training Pre-*  
28 *Publication* at 65.

- 1 • acquired through torrenting and direct downloading the Chicken Soup Works from
- 2 shadow-library repositories and datasets containing pirated works from shadow
- 3 libraries;
- 4 • distributed the Chicken Soup Works through the use of torrenting software,
- 5 programs, or protocols;
- 6 • reproduced additional copies during ingestion, preprocessing, storage, deduplication,
- 7 formatting, and/or tokenization; and
- 8 • while training its models made even more copies of the text—because every training
- 9 pass (each epoch and each step of gradient descent) automatically requires creating
- 10 and working with fresh versions of that text.

11 87. Defendant’s reproductions and distributions of Plaintiff’s copyrighted works were  
12 made without permission, license, or consent and violated Plaintiff’s exclusive rights under the  
13 Copyright Act.

14 88. Defendant’s infringement was **willful**. As alleged above, Defendant knowingly  
15 trained its models on and/or optimized its product with datasets saturated with pirated books,  
16 including the Chicken Soup Works; relied on shadow-library corpora it knew to be illegal; ignored  
17 internal and external warnings; attempted to conceal the composition of its training datasets; and  
18 continued copying after public reports, lawsuits, law-enforcement seizures, cease-and-desist  
19 notices, and industry-wide alerts made the illegality unmistakable.

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

22 90. Plaintiff has been injured by Defendant’s willful acts of copyright infringement.  
23 Plaintiff is entitled to statutory damages, actual damages, restitution of profits, and/or other remedies  
24 in law or equity.

25 91. Plaintiff is entitled to recover attorneys’ fees and costs under 17 U.S.C. § 505.

26 **COUNT II**

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

28 92. Plaintiff incorporates the allegations above.

1 93. Defendant used torrenting software, programs, or protocols to download datasets  
2 containing pirated copies of works, including the Chicken Soup Works.

3 94. In connection with its torrenting of datasets that contained copyrighted works,  
4 Defendant uploaded and distributed, either through “seeding” and/or “leeching,” copyrighted  
5 materials, including the Chicken Soup Works, thereby making those works available to third parties  
6 for downloading on peer-to-peer networks.

7 95. Defendant knowingly participated in peer-to-peer sharing networks that it knew  
8 trafficked in pirated copies of copyrighted materials. In other words, Defendant knew that others on  
9 these networks were infringing copyrighted materials through reproduction and/or distribution.  
10 There was no substantial or commercially significant non-infringing use of the copyrighted  
11 materials that Anthropic uploaded and distributed. Nor was there substantial or commercially  
12 significant non-infringing use of Defendant’s uploading and distribution of Plaintiff’s copyrighted  
13 works. By participating in these networks, and by further uploading and distributing Plaintiff’s  
14 copyrighted works, Defendant materially contributed to and induced further infringement of  
15 Plaintiff’s works.

16 96. By knowingly inducing and materially contributing to others’ infringement of  
17 Plaintiff’s works, Defendant is liable for contributory copyright infringement.

18 97. As a direct and proximate cause of Defendant’s conduct, Plaintiff was injured and is  
19 entitled to statutory damages, actual damages, restitution of profits, and/or other remedies in law or  
20 equity.

21 98. Plaintiff is entitled to recover attorneys’ fees and costs under 17 U.S.C. § 505.

22 **COUNT III**

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

24 99. Plaintiff incorporates the allegations above.

25 100. Plaintiff’s materials contain information that constitutes “copyright management  
26 information” as that term is defined in 17 U.S.C. § 1202(c). This includes but is not limited to author  
27 information, information about the copyright owner, and copyright notices.

28

1 101. Upon downloading copyrighted materials, including Plaintiff’s works, Defendant  
2 processed the data, and in doing so, removed and altered certain text and information, including  
3 copyright management information found in and on Plaintiff’s works. When it removed this  
4 information, Defendant did so without the authority of the copyright owners.

5 102. Defendant’s removal of the copyright management information was intentional—it  
6 did so to, among other things, create high-quality LLM training data and, through the creation and  
7 use of high-quality training data, ultimately create high-quality LLM models.

8 103. Defendant removed copyright management information from Plaintiff’s works  
9 knowing or having reasonable grounds to believe that it was enabling, facilitating, and concealing  
10 acts of copyright infringement. As to concealment, Defendant knew or had reasonable grounds to  
11 believe that by stripping copyrighted works of copyright management information it would be  
12 harder for others to discover the true sources—*e.g.*, copyrighted works—of Defendant’s training  
13 data.

14 104. Plaintiff was harmed by Defendant’s removal of copyright management information  
15 from its works and is entitled to statutory damages, actual damages, restitution of profits, and other  
16 remedies provided by law. Plaintiff is entitled to recover attorneys’ fees and costs under 17 U.S.C.  
17 § 1203(b)(5).

18 **PRAYER FOR RELIEF**

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

- 20 a. Judgment in favor of Plaintiff against the Defendant;
- 21 b. A declaration that the Defendant has infringed Plaintiff’s exclusive copyrights  
22 under the Copyright Act;
- 23 c. A declaration that such infringement is willful;
- 24 d. A declaration that Defendant violated 17 U.S.C. § 1202(b) through its removal  
25 of copyright management information;
- 26 e. A permanent injunction enjoining the Defendant and all those acting in concert  
27 with it from engaging in the infringing conduct alleged herein;

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

- f. That the Defendant be directed to account to Plaintiff for all gains, profits, and advantages derived from its unlawful acts;
- g. An award of statutory damages under the Copyright Act;
- h. An award of statutory or actual damages under 17 U.S.C. § 1203(c);
- i. An award of restitution, disgorgement, costs, expenses, and attorneys’ fees as permitted by law (including those allowable under 17 U.S.C. § 505 and/or 17 U.S.C. § 1203(b)(4)–(5));
- j. Pre- and post-judgment interest on the damages awarded to Plaintiff; and
- k. Further relief for Plaintiff as the Court may deem just and proper.

**JURY TRIAL DEMANDED**

Under Federal Rule of Civil Procedure 38(b), Plaintiff demands a trial by jury.

1 Dated: May 7, 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  
kroche@fnf.law  
vel@fnf.law  
apotter@fnf.law

*Counsel for Plaintiff*