

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. 44th Street, Ste 915
New York, NY 10017
T: (646) 494-2900

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

COGNELLA, INC.,

Plaintiff,

v.

ANTHROPIC PBC,

Defendant.

Civil Case No.:

COMPLAINT

DEMAND FOR JURY TRIAL

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

3 **I. INTRODUCTION**

4 2. This case concerns Defendant Anthropic’s exploitation of Plaintiff Cognella, Inc.’s
5 copyrighted works to build its “Claude” family of large language models (“LLMs”). Defendant is
6 in 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, Anthropic 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. Anthropic could
14 have—and *should have*—paid copyright owners, including Plaintiff, for licenses to use their
15 copyrighted works in connection with training its Claude models. Instead, Anthropic fed from, and
16 then back into, the illegal market for digitally available copyrighted materials. Ultimately, Anthropic
17 downloaded, reproduced, distributed, and defaced these copyrighted works, including Plaintiff’s
18 works, to get ahead in the generative AI arms race and improve its bottom line. Plaintiff brings this
19 action to hold Anthropic accountable for its brazen acts of copyright infringement.¹

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 ¹ 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. It’s 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 2026.² 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 ² 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.”³ 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 Cognella, Inc., an academic publisher, is a California stock corporation and
17 the owner or exclusive licensee of thousands of copyrights in textbooks, course readers, custom
18 course packs, and other learning materials.

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

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 ³ *Feltner v. Columbia Pictures Television, Inc.*, 523 U.S. 340, 353 (1998).

27 ⁴ 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 Cognella Works, the
16 reproduction and ingestion of those copies into Defendant’s training pipelines, the training and fine-
17 tuning of the relevant LLMs, and the commercialization of the resulting models—occurred in this
18 District.

19 **IV. FACTUAL ALLEGATIONS**

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

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

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

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

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

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

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

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

27
28

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

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

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

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

20 31. LLMs typically are trained, at least initially, by feeding them massive amounts of
21 text data from which they can “learn” statistical relationships between and among text. In the process
22 of compiling training data and using it during the training process, LLM developers as a matter of
23 course typically create multiple copies of “raw” datasets. So, for instance, if an LLM developer
24 downloads a copy of a digital book from some Internet-based source (*e.g.*, a shadow library), the

25
26 _____
27 ⁵ U.S. Copyright Office, *Copyright and Artificial Intelligence, Part 3: Generative AI Training Pre-*
Publication, at 17 (May 2025), <https://perma.cc/EY5U-EFUY>.

28 ⁶ *Id.*

1 developer typically will create a new “cleaned” copy of the book (*e.g.*, stripped of certain
2 undesirable data), deduplicate that book against other potential copies, compile that book with other
3 data to create a new compilation (or compilations), and store that book in multiple locations. In
4 short, with respect to training data, the LLM development pipeline typically involves numerous
5 reproductions of any given piece of the dataset.

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

10 33. In determining the corpus of training data for an LLM, a model’s developer typically
11 considers “the quantity of data, its quality, and the ultimate purpose(s) of the model.”⁸

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

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

23
24 ⁷ James Betker, *The “it” in AI Models is the Dataset*, (June 10, 2023), <https://perma.cc/ZCH9-H53S>.

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

27 ⁹ *Id.* at 10.

28 ¹⁰ *Id.* at 11.

1 36. Among other sources, academic content is extremely valuable as training material
2 for LLMs. This is evidenced by the fact that a number of the leading AI companies have been willing
3 to pay for it from other publishers.¹¹ For example, Microsoft recently paid Informa, the parent
4 company of Taylor & Francis, an initial fee of \$10 million to make use of its content “to help
5 improve relevance and performance of AI systems.”¹² Another leading academic publisher, Wiley,
6 recently agreed to sell academic content for training AI models by completing a “GenAI content
7 rights project” with an undisclosed “large tech company” in 2024.¹³ Wiley indicated that
8 negotiations for another generative AI project with a second large tech company were contemplated
9 for 2025.¹⁴

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

13 38. As discussed in the introduction, technology companies have treated generative AI
14 as the next foundational layer of the digital economy. Industry leaders publicly describe an “AI arms
15 race,” in which they have redirected their corporate strategies to seize control of what they believe
16 will become a new infrastructure layer for commerce, communication, and knowledge work.¹⁵ For
17

18 ¹¹ Roger C. Schonfeld, *Tracking the Licensing of Scholarly Content to LLMs*, THE SCHOLARLY
19 KITCHEN (Oct. 15, 2024) <https://perma.cc/9DX9-QEWW> (“A number of companies are in the hunt
20 for this content, including not only OpenAI and Google, but also Apple and more specialized
21 providers. Investment has been pouring in as a result of the market’s spike in interest in artificial
intelligence.”)

22 ¹² Christa Dutton, *Two Major Academic Publishers Signed Deals With AI Companies. Some
Professors Are Outraged*, THE CHRONICLE OF HIGHER EDUCATION, (July 29, 2024)
23 <https://perma.cc/PQV5-FWRB>.

24 ¹³ *Id.*; Press Release, WILEY, *Wiley Increases Quarterly Dividend for the 31st Consecutive Year*, (June
25 27, 2024) <https://perma.cc/JS5D-3WFN>.

26 ¹⁴ Dutton, *Two Major Academic Publishers Signed Deals With AI Companies. Some Professors Are
Outraged*.

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

1 these companies, staying ahead of competitors is “code red.”¹⁶ Among other things, being too slow
2 out of the blocks could mean ending up in last place; if by the time you publicly released *your* LLM
3 model, the public writ large and private industry had already adopted *other companies’* models,
4 there would be a real risk that your model would be left on the shelf.

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

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

9 40. For years now, there have been illicit, online marketplaces for digital copies of books.
10 These “shadow libraries,” as they are commonly known, systematically acquire, store, index, and
11 disseminate full-fidelity digital copies of copyrighted books—typically in native formats such as
12 EPUB, PDF, MOBI, or DJVU. These online repositories maintain searchable catalogs, metadata,
13 mirrors, bulk-download mechanisms, and—critically—complete downloadable archives designed
14 to allow third parties to replicate the entire collection.¹⁷ And they do all of this without authorization
15 from authors or publishers.¹⁸ In short, shadow libraries facilitate the reproduction and distribution
16 of unauthorized copies of copyrighted works at industrial scale.

17 41. These libraries are widely known within both piracy communities and the technology
18 sector as illegal sources of copyrighted books. Many have been the subject of criminal prosecutions,
19 civil injunctions, domain seizures, and formal designation as “notorious piracy markets” by United

21 ¹⁶ 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 ¹⁷ See Letter from Mary E. Rasenberger, CEO, Authors Guild, & Umair Kazi, Dir. of Pol’y &
24 Advocacy, Authors Guild, to Daniel Lee, Assistant U.S. Trade Representative for Innovation &
25 Intellectual Prop., Office of the U.S. Trade Representative (Oct. 7, 2022), <https://perma.cc/XM4R-NDN3>.

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

1 States trade authorities. As centralized shadow libraries increasingly faced enforcement actions,
2 including the seizure of domains, third parties responded by creating full mirrored copies of those
3 repositories for decentralized redistribution.¹⁹

4 42. In 2018, an OpenAI employee downloaded pirated copies of books from Library
5 Genesis, or “LibGen”—a shadow library repeatedly enjoined by federal courts²⁰—and used those
6 books to create two internal datasets OpenAI called “LibGen1” and “LibGen2,” which OpenAI
7 publicly referred to as “Books1” and “Books2.”²¹ OpenAI used those pirated corpora to train GPT-
8 3, which it released in June 2020 to widespread commercial acclaim. OpenAI’s use of a books
9 corpora from a pirate library established the model for how to quickly and cheaply obtain published
10 materials for use as LLM training data.

11 43. Within weeks of GPT-3’s release, an open research collective, EleutherAI, formed
12 with the express goal of replicating GPT-3’s capabilities and democratizing access to large-scale
13 language modeling. To do so, EleutherAI assembled and publicly released *The Pile*, a large (800+
14 gigabyte) general-purpose training dataset expressly designed to be downloaded, used, and
15 incorporated into LLMs by academic researchers, startups, and commercial AI developers.²²

16 44. Because OpenAI did not disclose the precise makeup of its training datasets,
17 members of EleutherAI constructed a pirated book corpus of their own: “Books3,” consisting of
18
19
20
21

22 ¹⁹ Woodcock, *‘Shadow Libraries’ Are Moving Their Pirated Books to the Dark Web After Fed Crackdowns*.

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

25 ²¹ *In re OpenAI, Inc., Copyright Infringement Litig.*, Case No. 1:2025-md-03143 (S.D.N.Y. 2025),
26 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>.

27 ²² Leo Gao et al., *The Pile: An 800GB Dataset of Diverse Text for Language Modeling*, arXiv, 1
28 (2020), available at <https://perma.cc/NHV6-R8YE>.

1 approximately 196,640 books.²³ Books3 comprises about 12 percent (just over 100 gigabytes) of
2 *The Pile*.²⁴

3 45. EleutherAI and the compiler of Books3, Shawn Presser, have confirmed the genesis
4 of Books3 is the shadow library Bibliotik. Presser has publicly stated that Books3 represents “all of
5 bibliotik,”²⁵ and an EleutherAI paper likewise confirms that Books3 “is a dataset of books derived
6 from a copy of the contents of the Bibliotik private tracker.”²⁶

7 46. Bibliotik is a private, invitation-only torrent tracker that has long functioned as a
8 centralized source of pirated e-books. The repository, which hosts and distributes hundreds of
9 thousands of copyrighted books, is only accessible with the use of torrenting protocols, which are
10 explained below.²⁷

11 47. Bibliotik has been widely recognized in piracy communities, academic literature, and
12 AI-research documentation as a shadow library devoted to copyrighted books. Unsurprisingly, then,
13 *The Pile*’s datasheet acknowledges that “Books3 is almost entirely comprised of copyrighted
14 works.”²⁸ Bibliotik’s illicit nature is no secret—it has been openly discussed for years prior to major
15 tech companies’ use of datasets derived directly from it.²⁹

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

18

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

20 ²⁴ Biderman et al., *Datasheet for the Pile*, arXiv, 8.

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

22 ²⁶ Biderman et al., *Datasheet for the Pile*, arXiv, 8.

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

25 ²⁸ Biderman et al., *Datasheet for the Pile*, arXiv, 8.

26 ²⁹ Kyle Barr, *Anti-Piracy Group Takes Massive AI Training Dataset ‘Books3’ Offline*, GIZMODO
27 (Aug. 18, 2023, 8:50 AM), <https://perma.cc/5ZL9-RQCQ>; Peter Schoppert, *Whether you’re an
28 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>.

1 [Presser] created Books3 in the hope that it would allow any developer to create
2 generative-AI tools. “It would be better if it wasn’t necessary to have something like
3 Books3,” he said. “But the alternative is that, without Books3, only OpenAI can do
4 what they’re doing.”³⁰

5 49. EletheurAI’s compilation and distribution of *The Pile* and Books3 provided “off-the-
6 shelf” access to a corpus of infringing works that any AI developer could download and immediately
7 incorporate into a large-scale training pipeline. As a result, Books3’s presence within *The Pile*
8 facilitated wide downstream distribution and adoption of a corpus derived from a pirate book
9 library.³¹

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

12 51. As noted above, OpenAI’s developers trained its LLMs on a books corpus derived
13 from LibGen, one of the largest and longest-running shadow libraries in the world. LibGen hosts
14 millions of pirated books, academic texts, and scholarly articles.³² It operates as a centralized
15 repository offering direct downloads of full-fidelity e-book files. It also distributes its entire database
16 through bulk archives and torrent files, enabling third parties to download and locally host complete
17 copies of its collection.³³ LibGen has been repeatedly enjoined by federal courts for copyright
18 infringement and has been designated a “notorious market” by the United States Trade
19
20

21 ³⁰ *Id.*

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

24 ³² Office of the U.S. Trade Representative, REVIEW OF NOTORIOUS MARKETS FOR
25 COUNTERFEITING AND PIRACY, 27 (2024), <https://perma.cc/22VN-9VD7>. (“Libgen ... hosts a
26 large number of digital copies of books, manuals, journals, and other works, many of which are
unauthorized copies of copyright protected content.”)

27 ³³ See Letter from Mary E. Rasenberger, CEO, Authors Guild, & Umair Kazi, Dir. of Pol’y &
28 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.

1 Representative.³⁴ Despite enforcement actions, LibGen has remained accessible through shifting
2 domains, mirrors, and downloadable archives. Its persistence is a function of deliberate
3 decentralization designed to evade shutdown.³⁵

4 52. Z-Library (also known as “B-ok”) is another well-known shadow library that
5 emerged as an expanded and user-friendly derivative of LibGen. It incorporated large portions of
6 LibGen’s catalog while adding additional titles, metadata, and interface features.³⁶ Z-Library
7 offered premium features—including faster downloads and higher volume limits—in exchange for
8 payment, operating in effect as a commercial piracy service.³⁷ In 2022, Z-Library’s domains were
9 seized by law-enforcement authorities, and its operators were arrested and later indicted for criminal
10 copyright infringement.³⁸ These actions confirmed what had long been publicly known: Z-Library
11 was an illegal piracy operation. The seizure of Z-Library did not eliminate access to its content.
12 Instead, third parties responded by creating full mirrors of its collection to ensure continued
13 distribution.³⁹

14 53. Another major player in the shadow library ecosystem is PiLiMi, which is a complete
15 mirrored archive of the Z-Library corpus, explicitly created to preserve and propagate Z-Library’s
16
17
18

19 ³⁴ See Office of the U.S. Trade Representative, 2019 REVIEW OF NOTORIOUS MARKETS FOR
20 COUNTERFEITING AND PIRACY, 27, <https://perma.cc/22VN-9VD7>.

21 ³⁵ See Andrew Albanese, *Textbook Publishers Sue Notorious ‘Shadow Library’ Libgen*,
22 PUBLISHERS WEEKLY (Sep. 14, 2023), <https://perma.cc/3NPY-UJCX>.

23 ³⁶ Jordana Rosenfeld, *Z-Library*, ENCYCLOPEDIA BRITANNICA (last accessed March 9, 2026),
<https://perma.cc/4H26-GDCC>.

24 ³⁷ See Masood Farivar, *Two Russian Nationals Charged With Operating E-Book Piracy Site*, VOA
25 (Nov. 16, 2022), <https://perma.cc/6MPD-QNKB>.

26 ³⁸ Press Release, U.S. Dep’t of Justice, U.S. Att’y’s Off., E. Dist. of N.Y., *Two Russian Nationals*
27 *Charged with Running Massive E-Book Piracy Website* (Nov. 16, 2022), [https://perma.cc/CR4L-](https://perma.cc/CR4L-JLA3)
[JLA3](https://perma.cc/CR4L-JLA3).

28 ³⁹ See Woodcock, *‘Shadow Libraries’ Are Moving Their Pirated Books to The Dark Web*.

1 pirated collection after law-enforcement seizures, ensuring continuity of access despite shutdowns
2 of the original site.⁴⁰

3 54. PiLiMi is not merely a website or index. It is a full, downloadable dataset designed
4 to allow users to obtain and locally host millions of pirated books through torrent distribution.⁴¹
5 Users who download PiLiMi do not passively receive data; they actively participate in copying and
6 redistributing copyrighted works through torrent “leeching” and “seeding.”⁴²

7 55. PiLiMi later rebranded as “Anna’s Archive” and expanded to aggregate and host the
8 complete collections of LibGen, Z-Library, and other pirated sources.⁴³ Like PiLiMi, Anna’s
9 Archive functions as a meta-library: it indexes, mirrors, and redistributes multiple shadow libraries
10 simultaneously, offering users unified access to millions of pirated books.⁴⁴

11 56. Anna’s Archive offers paid tiers that provide “high-speed” or priority access to its
12 pirated collections.⁴⁵ Through its downloadable archives and torrent-based distribution, Anna’s
13 Archive enables users to acquire and store local copies of millions of copyrighted books in bulk.⁴⁶

14
15
16 ⁴⁰ See Ernesto Van de Sar, “*Anna’s Archive*” *Opens the Door to Z-Library and Other Pirate*
17 *Libraries*, TORRENTFREAK (Nov. 19, 2022), <https://perma.cc/U88R-WTR4>.

18 ⁴¹ See Geoff Wheelright, *Will I get a piece of Anthropic’s \$1.5B settlement if my book was used to*
19 *train AI?*, GEEKWIRE (Sep. 18, 2025), <https://perma.cc/X5EW-TKY3>.

20 ⁴² See Robert Nogacki, *Anthropic’s Landmark Settlement: A \$1.5 Billion Copyright Precedent in*
21 *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
22 users simultaneously download and distribute files”).

23 ⁴³ See Ernesto Van de Sar, “*Anna’s Archive*” *Opens the Door to Z-Library and Other Pirate*
Libraries.

24 ⁴⁴ *Id.*

25 ⁴⁵ See *If you’re an LLM, please read this*, ANNA’S ARCHIVE (February 18, 2026),
26 <https://perma.cc/989X-GS3Q>.

27 ⁴⁶ 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-MYBY>.

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

3 58. According to Anna’s Archive, “virtually all major companies building LLMs
4 contacted us to train on our data. . . . We have given high-speed access to about 30 companies.”⁴⁷
5 Anna’s Archive blog stated as recently as February 18, 2026 that if an *LLM was reading its blog*
6 “you have likely been trained in part on our data.”⁴⁸

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

19 **D. Anthropic’s Deliberate Infringement of Plaintiff’s Copyrights**

20 **1. Anthropic’s Business and Bypassing of the Licensing Market**

21 60. Anthropic’s business model is built on the large-scale copying of published content.
22 *First*, Anthropic developed and commercialized the Claude family of large language models by
23 stealing up to seven million copyrighted books, including the Cognella Works.⁴⁹ *Second*, Anthropic

24 _____
25 ⁴⁷ See *Copyright reform is necessary for national security*, ANNA’S ARCHIVE (Jan. 31, 2025),
<https://perma.cc/3RVZ-6G5B>.

26 ⁴⁸ See *If you’re an LLM, please read this*, ANNA’S ARCHIVE (Feb. 18, 2026), [https://perma.cc/MJ7Z-](https://perma.cc/MJ7Z-3ZCL)
27 [3ZCL](https://perma.cc/MJ7Z-3ZCL).

28 ⁴⁹ See *Bartz v. Anthropic PBC*, 791 F. Supp. 3d 1038, 1046 (N.D. Cal. 2025).

1 also had the explicit goal to “amass a central library of ‘all the books in the world’ to retain
2 ‘forever.’”⁵⁰

3 61. To accomplish its goal, “it stole the works for its central library by downloading them
4 from pirated libraries.”⁵¹ Anthropic did explore the licensing market but quickly decided that it was
5 not “a practical approach”⁵² to satisfy its thirst for copyrighted works. Judge Alsup, presiding over
6 the *Bartz v. Anthropic* litigation, stated simply: “From the start, Anthropic ‘ha[d] many places from
7 which’ it could have purchased books, but it preferred to steal them to avoid ‘legal/practice/business
8 slog.’”⁵³ In short, rather than paying for copyrighted materials, Anthropic downloaded pirated
9 copies of protected works, reproduced them, fed them into its models, and otherwise retained them
10 for its own “central library.”

11 2. Anthropic’s Acquisition of Pirated Works from Shadow Libraries

12 62. To train Claude, Anthropic downloaded Books3 in 2021, which co-founder Ben
13 Mann “knew had been assembled from unauthorized copies of copyrighted books,” downloaded at
14 least five million copies of books from LibGen, “which [Mann] knew had been pirated,” and
15 downloaded at least two million copies of books from PiLiMi, which “Anthropic knew had been
16 pirated.”⁵⁴

17 63. Anthropic’s decision to base its flagship models on pirated books was driven by
18 commercial advantage. As Jared Kaplan, Anthropic’s co-founder and Chief Science Officer, has
19 explained, “it is important to obtain vast amounts of books and also to have diverse types of books
20 in the training corpus to create a model with truly generative capabilities.”⁵⁵ As long-form content,

21 _____
22 ⁵⁰ *Bartz v. Anthropic PBC*, 787 F. Supp. 3d 1007, 1014 (N.D. Cal. 2025).

23 ⁵¹ *Id.* at 1029.

24 ⁵² *See Bartz v. Anthropic PBC*, No. 3:24-cv-05417 (N.D. Cal. June 23, 2025), ECF No. 553, at 17-
25 18.

26 ⁵³ *Bartz v. Anthropic PBC*, 787 F. Supp. 3d at 1015.

27 ⁵⁴ *Id.*

28 ⁵⁵ *See* Kaplan Decl. ¶ 47, *Bartz*, ECF No. 128.

1 training LLMs on the “entire text” of published material—as Anthropic has admitted to doing—
2 offers great value.⁵⁶

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

5 3. Anthropic’s Torrenting and Distribution of Pirated Works

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

8 66. In June 2021, Mann personally torrented approximately five million pirated books
9 from LibGen for Anthropic’s use.⁵⁷ Mann acted with the knowledge and approval of Anthropic’s
10 senior leadership. Before the LibGen torrenting, he discussed the plan with co-founders Dario
11 Amodei, Jared Kaplan, and other senior leaders.

12 67. Anthropic’s own Archive Team had described LibGen as a “blatant violation of
13 copyright,” and Amodei himself called it “sketchy.”⁵⁸ Yet Anthropic proceeded, choosing torrenting
14 over purchasing or licensing the copyrighted materials.

15 68. Anthropic repeated the same conduct in 2022 with PiLiMi. As U.S. law enforcement
16 was working to shut down existing pirate libraries, a group online copied LibGen and built upon it
17 to create Z-Library. The FBI later shut down Z-Library as well.⁵⁹ However, by that point, Z-Library
18 had itself been fully copied, or “mirrored,” into another repository: PiLiMi.⁶⁰ Mann circulated the
19 PiLiMi source to colleagues, and Anthropic employees torrented approximately two million
20
21
22

23 ⁵⁶ *Id.* at ¶¶ 43, 47.

24 ⁵⁷ B. Mann Dep. Tr. at 89:6-8, (Aug. 15 and 18, 2025), *Bartz*, ECF No. 337-1.

25 ⁵⁸ *Id.* at 144:4-13, 396:3-13.

26 ⁵⁹ *See Bartz*, 791 F. Supp. 3d at 1046.

27 ⁶⁰ *Id.*

1 additional pirated books not already captured from LibGen.⁶¹ Those files were not abstract “data”
2 but full-text digital books in .epub, .pdf, and .txt formats,⁶² including the Cognella Works.

3 69. In torrenting from these files, Anthropic knew or should have known that it was
4 participating in a *peer-to-peer network that traded in unauthorized copies of copyrighted material*.
5 Stated differently, Anthropic knew or should have known that it was facilitating further copyright
6 infringement by making available to others on the peer-to-peer network unauthorized copies of
7 copyrighted materials.

8 4. Anthropic’s Stripping of Copyright Management Information

9 70. Anthropic not only reproduced and distributed the Cognella Works without
10 authorization, it also deliberately stripped those materials of copyright management information
11 (“CMI”) before using them to train Claude.

12 71. Anthropic knew that major AI-training datasets, including *The Pile*, WebText,
13 WebText2, and Common Crawl, contained copyrighted works whose copyright notices, ownership
14 information, and other identifying material had been removed through extraction tools such as
15 Newspaper, Dragnet, Readability, and jusText. Anthropic’s founders and senior employees were
16 familiar with those tools before and after Anthropic’s founding. Dario Amodei, Benjamin Mann,
17 Jared Kaplan, and other future Anthropic personnel had used or developed datasets at OpenAI that
18 extracted text from scraped webpages while omitting surrounding material, including footers where
19 copyright notices typically appear.⁶³ Yet Anthropic trained Claude on these same stripped datasets,
20 despite knowing that they included unauthorized copies of copyrighted works stripped of their CMI.

21 72. Anthropic also affirmatively chose tools that removed CMI more effectively,
22 including from Plaintiff’s copyrighted materials. In 2021, Mann, Kaplan, and other Anthropic
23 employees compared extraction tools for filtering training data and then used these extraction

24 ⁶¹ *Id.*

25 ⁶² *Id.*

26
27 ⁶³ 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 methods to clean its training data, including the Cognella Works. The purpose and effect of this
2 process was to train Claude on Plaintiff’s expressive content while suppressing the ownership
3 information attached to it. Anthropic wanted Claude to reproduce protected expression, not the
4 copyright notices that identify the rights holders. By stripping that information from training data
5 and outputs, Anthropic concealed the source of its infringement from users, Plaintiff, and other
6 copyright owners.

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

8 73. Since 2024, Anthropic has purchased physical books at scale, often in batches of tens
9 of thousands, and scanned them into digital files for AI training.⁶⁴ To date, Anthropic has spent
10 millions of dollars buying and scanning millions of physical books.⁶⁵ This scanning project was not
11 authorized by Plaintiff or other copyright owners. As with the millions of pirated books Anthropic
12 torrented from LibGen and PiLiMi, Anthropic converted these works into training material without
13 permission, payment, or license.

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

15 74. Defendant did not merely make temporary copies of copyrighted works for a discrete
16 training process (not that that would be permissible, in any event). Instead, from all those
17 infringements, Anthropic created a general “research library” or “generalized data area.”⁶⁶ That is,
18 Anthropic created a permanent central library of pirated published materials, including the Cognella
19 Works, to exploit for a range of unspecified *other* purposes, such as “research.”⁶⁷ Anthropic
20 continued this conduct, for which it “lacked any entitlement to hold copies of the books at all” and
21 “retain[ed] them even after deciding it would not make further copies from them for training.”⁶⁸ The
22 plan was clear: “ke[eping] in the original version of the underlying book files Anthropic had

23 _____
24 ⁶⁴ See *Bartz*, 791 F. Supp. 3d at 1047.

25 ⁶⁵ *Bartz*, ECF No. 553, at 17-18.

26 ⁶⁶ *Bartz*, 787 F. Supp. 3d at 1016.

27 ⁶⁷ *Id.*

28 ⁶⁸ *Id.* at 1031.

1 obtained or created, that is, pirated or scanned” and “stor[ing] everything forever.”⁶⁹ Anthropic saw
2 “no compelling reason to delete a book—even if not used for training LLMs.”⁷⁰

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

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

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

12 77. A 2025 study by researchers from Cornell, Stanford, and West Virginia University
13 tested leading LLM models, including Anthropic’s Claude 3.7 Sonnet, and extracted memorized
14 copyrighted books from them.⁷¹ Using repeated prompting and continuation techniques, the
15 researchers were able to recover lengthy, near-verbatim passages. The results for Claude were
16 severe. From Claude 3.7 Sonnet alone, researchers extracted 97.5% of *The Great Gatsby*, 95.5% of
17 *1984*, 94.3% of *Frankenstein*, 92.3% of *Harry Potter and the Sorcerer’s Stone*, and 70.2% of *The*
18 *Hobbit*.⁷² These results show that Claude does not merely learn abstract patterns from copyrighted
19 materials; it stores recoverable copies of the works at extraordinary scale.

20 78. A separate 2026 study reached the same conclusion through a different method.
21 Researchers from Stony Brook University, Carnegie Mellon University, and Columbia Law School
22 found that finetuned frontier models could reproduce up to 85–90% of copyrighted books, with
23

24 ⁶⁹ *Id.* at 1016 (internal quotation marks omitted).

25 ⁷⁰ *Id.* (internal quotation marks omitted).

26 ⁷¹ A. Feder Cooper et al., *Extracting Copyrighted Long-Form Text from Production Language*
27 *Models*, arXiv, 1-2 (2025), <https://perma.cc/F9YQ-NKV7>.

28 ⁷² *Id.* at 12, Fig. 5.

1 single verbatim spans exceeding 460 words.⁷³ The models reproduced these passages from semantic
2 prompts alone, without being given the book text in the prompt.

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

9 80. Together, these studies also show that safety filters do not remove the underlying
10 copyrighted material from the model. The works remain embedded in model weights and can be
11 extracted despite protective layers.⁷⁴ In practical effect, Claude functions as a repository of pirated
12 copyrighted works: Anthropic copied published materials to train the model, and the model retained
13 near-verbatim reproductions of them.

14 **8. Anthropic's Harm to the Market for Plaintiff's Copyrighted Materials**

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

21 82. *First*, Anthropic's decision to download and use unauthorized copies of the Cognella
22 Works from shadow libraries deprived Plaintiff of revenue in the form of licensing fees that it would
23 have otherwise earned. As alleged herein, there is an existing market for licensing copyrighted
24 materials such as Plaintiff's, including for use in the development of LLMs. Anthropic bypassed
25

26 ⁷³ A. Liu et al., *Alignment Whack-a-Mole : Finetuning Activates Verbatim Recall of Copyrighted*
27 *Books in Large Language Models*, arXiv, 2, (2026), <https://perma.cc/EC6X-Z38M>.

28 ⁷⁴ *Id.* at 9–11.

1 that market, and in doing so, deprived Plaintiff of licensing revenue it would have earned. Anthropic
2 now allows users to opt out of having their data used to train its AI models, but it deprived Plaintiff
3 and many others of that choice. Anthropic’s misconduct undermined Plaintiff and many others,
4 eliminating the bargaining power they should have had, and otherwise would have had, with respect
5 to licensing terms for the use Anthropic made of their works.

6 83. *Second*, Claude, trained on the Cognella Works (and the protected works of others),
7 is capable of generating outputs that compete directly with, and risk serving as replacements for, the
8 Cognella Works. As the U.S. Copyright Office has warned, “the speed and scale at which AI systems
9 generate content pose a serious risk of diluting markets for works of the same kind as in their training
10 data.”⁷⁵

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

15 **V. CLAIMS FOR RELIEF**

16 **COUNT I**

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

18 85. Plaintiff incorporates the allegations above.

19 86. Plaintiff is the legal or beneficial owner of the copyrighted works listed in Exhibit A
20 (referred to herein as the Cognella Works).

21 87. The Defendant, without authorization from Plaintiff, copied, downloaded,
22 reproduced, ingested, parsed, embedded, and used pirated copies of the Cognella Works in the
23 development, training, fine-tuning, and deployment of their commercial large language models.
24 These acts violated Plaintiff’s exclusive rights under 17 U.S.C. § 106.
25

26
27 ⁷⁵ US Copyright Office, *Copyright and Artificial Intelligence, Part 3: Generative AI Training Pre-*
28 *Publication* at 65.

1 88. Defendant’s infringement occurred repeatedly throughout the lifecycle of its AI-
2 model development. As alleged above, Defendant:

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

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

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

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

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

27 93. Plaintiff is entitled to recover attorneys’ fees and costs under 17 U.S.C. § 505.
28

1 **COUNT II**

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

3 94. Plaintiff incorporates the allegations above.

4 95. Defendant used torrenting software, programs, or protocols to download datasets
5 containing pirated copies of works, including the Cognella Works.

6 96. In connection with its torrenting of datasets that contained copyrighted works,
7 Defendant uploaded and distributed, either through “seeding” and/or “leeching,” copyrighted
8 materials, including the Cognella Works, thereby making those works available to third parties for
9 downloading on peer-to-peer networks.

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

19 98. By knowingly inducing and materially contributing to others’ infringement of
20 Plaintiff’s works, Defendant is liable for contributory copyright infringement.

21 99. As a direct and proximate cause of Defendant’s conduct, Plaintiff was injured and is
22 entitled to statutory damages, actual damages, restitution of profits, and/or other remedies in law or
23 equity.

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

25 **COUNT III**

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

27 101. Plaintiff incorporates the allegations above.

28

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

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

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

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

17 106. Plaintiff was harmed by Defendant’s removal of copyright management information
18 from its works and is entitled to statutory damages, actual damages, restitution of profits, and other
19 remedies provided by law. Plaintiff is entitled to recover attorneys’ fees and costs under 17 U.S.C.
20 § 1203(b)(5).

21 **PRAYER FOR RELIEF**

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

- 23 a. Judgment in favor of Plaintiff against the Defendant;
- 24 b. A declaration that the Defendant has infringed Plaintiff’s exclusive copyrights
25 under the Copyright Act;
- 26 c. A declaration that such infringement is willful;
- 27 d. A declaration that Defendant violated 17 U.S.C. § 1202(b) through its removal
28 of copyright management information;

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

- e. A permanent injunction enjoining the Defendant and all those acting in concert with it from engaging in the infringing conduct alleged herein;
- f. That the Defendant be directed to account to Plaintiff for all gains, profits, and advantages derived from their 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 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. 44th Street, Ste 915
New York NY 10017
T: (646) 494-2900
vel@fnf.law
kroche@fnf.law
apotter@fnf.law

Counsel for Plaintiff