

1 William Dreher (*pro hac vice*)  
 2 Derek W. Loeser (*pro hac vice*)  
 3 Cari C. Laufenberg (*pro hac vice*)  
 4 Benjamin Gould (SBN 250630)  
 KELLER ROHRBACK L.L.P.  
 1201 Third Avenue, Suite 3400  
 Seattle, WA 98101-3268  
 5 (206) 623-1900  
 wdreher@kellerrohrback.com  
 6 dloeser@kellerrohrback.com  
 7 claufenberg@kellerrohrback.com  
 bgould@kellerrohrback.com  
 8  
 Matthew Butterick (SBN 250953)  
 9 BUTTERICK LAW PC  
 1920 Hillhurst Avenue, #406  
 10 Los Angeles, CA 90027  
 11 (323) 968-2632  
 mb@buttericklaw.com

Rohit D. Nath (SBN 316062)  
 SUSMAN GODFREY L.L.P.  
 1900 Avenue of the Stars, Suite 1400  
 Los Angeles, CA 90067  
 (310) 789-3100  
 RNath@susmangodfrey.com

Justin A. Nelson (*pro hac vice*)  
 Alejandra C. Salinas (*pro hac vice*)  
 SUSMAN GODFREY L.L.P.  
 1000 Louisiana Street, Suite 5100  
 Houston, TX 77002  
 (713) 653-7802  
 JNelson@susmangodfrey.com  
 ASalinas@susmangodfrey.com

***Additional Counsel Listed in Signature Block***

UNITED STATES DISTRICT COURT  
 NORTHERN DISTRICT OF CALIFORNIA  
 OAKLAND DIVISION

GRADY HENDRIX ET AL.,

*Individual and Representative Plaintiffs,*

v.

APPLE INC.,

*Defendant.*

No. 4:25-cv-07558-YGR

**CONSOLIDATED CLASS ACTION  
 COMPLAINT**

**JURY TRIAL DEMANDED**

This Document Relates To:  
 All Consolidated Actions

Plaintiffs Grady Hendrix, Jennifer Roberson, John Hornor Jacobs, and Eboni McKinnon (a/k/a Kianna Alexander) (together “Plaintiffs”), on behalf of themselves and all others similarly situated (the “Class,” as defined below), bring this class-action complaint (“Complaint”) against Defendant Apple Inc. (“Apple” or “Defendant”).

## I. INTRODUCTION

1  
2 1. Apple, one of the world’s most valuable companies, has invested substantial capital and  
3 engineering resources into developing generative artificial intelligence (“AI”) programs and  
4 technologies. Among Apple’s AI products is Apple Intelligence. Apple regards Apple Intelligence as a  
5 breakthrough innovation that will make its users’ experiences “profoundly different” across various  
6 product applications. Through Apple Intelligence, Apple hopes to add trillions to its market  
7 capitalization in coming years.

8 2. But Apple is building this new enterprise in part by using Books3, a dataset of pirated  
9 copyrighted books that includes the published works of Plaintiffs and the Class. Apple has developed,  
10 designed, maintained, and commercialized multiple large language models (“LLMs”). An LLM is a  
11 type of machine-learning model designed to emit convincingly naturalistic text outputs in response to  
12 user prompts. Companies developing LLMs “train” the LLMs to produce these outputs by feeding the  
13 models copies of an enormous quantity of data—frequently including massive numbers of textual  
14 works. Apple used the dataset called Books3 to train one set of LLMs, the OpenELM language models.  
15 Apple also, on information and belief, trained a second set of LLMs, its Foundation Language Models,  
16 using this same pirated dataset.

17 3. Apple’s generative AI also uses Applebot, a software program that copies mass  
18 quantities of webpages (also known as “scraping”). Apple scraped data with Applebot for nearly nine  
19 years before disclosing that it intended to train its commercial AI products on this scraped data.  
20 Scrapers like Applebot can also reach so-called “shadow libraries” that host millions of other  
21 unlicensed copyrighted books, including, on information and belief, Plaintiffs’ and Class Members’  
22 copyrighted works. Shadow libraries are the approximate equivalent of Napster for books; users upload  
23 unauthorized copies of books, and the shadow library helps other users download and access the works  
24 for free without paying the customary price. As has become public in other cases, a number of Apple’s  
25 competitors in the AI race—including OpenAI, Meta Platforms, Anthropic, and NVIDIA—resorted to  
26 obtaining copyrighted materials from shadow libraries like Library Genesis, Z-Library, Anna’s  
27 Archive, Sci-Hub, and Pirate Library Mirror. Upon information and belief, like its competitors, Apple  
28 obtained works from these notorious pirate libraries—accessed either through Applebot or torrenting—

1 to build its own permanent library of stolen material for training its models, for commercial research,  
2 and for other purposes.

3 4. The quality of Apple’s models depends on the data on which it is trained. LLMs operate  
4 by copying and later simulating creative expression found in copyrighted works. For this reason, the  
5 inclusion of expressive, high-quality material—especially copyrighted material—in Apple’s AI training  
6 datasets is deliberate and commercially significant. For instance, to access even more copyrighted  
7 material to develop its valuable generative AI products, Apple entered into a multimillion-dollar  
8 licensing agreement with Shutterstock. But not with Plaintiffs or the Class.

9 5. Plaintiffs and the Class are owners of registered copyrights for their published works.  
10 They did not consent to the use of their works in any model developed by Apple, including the  
11 Foundation Language Models and OpenELM language models.

12 6. The licensing market for AI training data is robust and growing. For instance, in  
13 February 2026, Microsoft launched a data-licensing marketplace called the Publisher Content  
14 Marketplace; Amazon has held discussions with publishers about building a similar service.  
15 Nevertheless, Apple did not compensate creators for use of their copyrighted works and concealed the  
16 sources of the datasets they downloaded, curated, and used for AI training to evade legal scrutiny. On  
17 information and belief, Apple continues to retain its own curated, private shadow library of pirated  
18 books to train AI models currently or soon to be under development and for other commercial  
19 purposes, all without seeking Plaintiffs’ or Class Members’ consent or providing them compensation.

20 7. Apple has copied the copyrighted works of Plaintiffs and the Class to train AI models  
21 whose outputs compete with and dilute the market for those very works. Without those works, Apple’s  
22 generative AI models, Apple products that integrate those models via Apple Intelligence, and Apple  
23 ecosystems in which third-party products integrate those models (such as the Apple App Store) would  
24 have far less commercial value. This conduct has deprived Plaintiffs and the Class of control over their  
25 work, undermined the economic value of their labor, and positioned Apple to achieve massive  
26 commercial success through unlawful means.

1 **II. JURISDICTION AND VENUE**

2 8. This Court has subject-matter jurisdiction under 28 U.S.C. § 1331 because this case  
3 arises under the Copyright Act (17 U.S.C. § 501).

4 9. Jurisdiction and venue are proper in this judicial district under 28 U.S.C. § 1391(c)(2)  
5 because Apple is headquartered in this district.

6  
7 **III. DIVISIONAL ASSIGNMENT**

8 10. Under Civil Local Rule 3-2(c), assignment of this case to the Oakland Division is  
9 proper because this case pertains to intellectual-property rights, which is a district-wide case category  
10 under General Order No. 44, and therefore venue is proper in any courthouse in this District.

11  
12 **IV. PARTIES**

13 **A. Plaintiffs**

14 11. Plaintiff Grady Hendrix is an author who lives in New York. He owns registered  
15 copyrights in the following works that are within the datasets that Apple reproduced without  
16 permission (collectively, the “Hendrix Asserted Works”): *Horrorstor* (U.S. Copyright Office  
17 Registration Number TX0007957438); *My Best Friend’s Exorcism* (USCO Registration Number  
18 TX0008276604); *Paperbacks from Hell* (USCO Registration Number TX0008461565); and *We Sold*  
19 *Our Souls* (USCO Registration Number TX0008645012).

20 12. Plaintiff Jennifer Roberson is an author who lives in Arizona. She owns registered  
21 copyrights in the following works that are within the datasets that Apple reproduced without  
22 permission (collectively, the “Roberson Asserted Works”): *Life and Limb* (USCO Registration Number  
23 TX0008936965); *Shapechangers* (USCO Registration Number TX0001405603); *Legacy of the Sword*  
24 (USCO Registration Number TX0002124929); *Track of the White Wolf* (USCO Registration Number  
25 TX0002138857); *A Pride of Princes* (USCO Registration Number TX0002595154); *Daughter of the*  
26 *Lion* (USCO Registration Number TX0002655076); *Flight of the Raven* (USCO Registration Number  
27 TX0003038048); *A Tapestry of Lions* (USCO Registration Number TX0003445261); *Karavans*  
28 (USCO Registration Number TX0006494921); *Deepwood* (USCO Registration Number

1 TX0006608590); *The Wild Road* (USCO Registration Number TX0007614823); *Lady of the Glen*  
2 (USCO Registration Number TX0004292057); *Lady of the Forest* (USCO Registration Number  
3 TX0003467465); *Lady of Sherwood* (USCO Registration Number TX0005700486); *Sword-Dancer*  
4 (USCO Registration Number TX0002125935); *Sword-Singer* (USCO Registration Number  
5 TX0002585643); *Sword-Maker* (USCO Registration Number TX0003031154); *Sword-Breaker* (USCO  
6 Registration Number TX0003187269); *Sword-Born* (USCO Registration Number TX0004752101);  
7 *Sword-Sworn* (USCO Registration Number TX0005504873); *Sword-Bound* (USCO Registration  
8 Number TX0007724532); *The Song of Homana* (USCO Registration Number TX0001840915).

9 13. Plaintiff John Hornor Jacobs is an author who lives in Arkansas. He owns registered  
10 copyrights in the following works that are within the datasets that Apple reproduced without  
11 permission (collectively, the “Hornor Jacobs Asserted Works”): *The Conformity* (USCO Registration  
12 Number TX0008114163); *The Sea Dreams It Is the Sky* (USCO Registration Number TX0008830088);  
13 *The Shibboleth* (USCO Registration Number TX0007905434); and *This Dark Earth* (USCO  
14 Registration Number TX0007561564).

15 14. Plaintiff Eboni McKinnon (a/k/a Kianna Alexander) is an author who lives in North  
16 Carolina. She owns registered copyrights in the following works that are within the datasets that Apple  
17 reproduced without permission (collectively, the “McKinnon Asserted Works”): *Love for All Time*  
18 (USCO Registration Number TX0008672549); *Couldn't Ask for More* (USCO Registration Number  
19 TX0008668191); and *Tempo of Love* (USCO Registration Number TX0008448578).

20 15. Below, the Hendrix Asserted Works, Roberson Asserted Works, Hornor Jacobs Asserted  
21 Works, and McKinnon Asserted Works are referred to collectively as the “Infringed Works.” On  
22 information and belief, Apple may have also infringed other works by Plaintiffs. A complete list of all  
23 of the Plaintiffs’ registered copyrighted works is attached as Exhibit A to this Complaint.

## 24 **B. Defendant**

25 16. Defendant Apple Inc. is a California corporation with its principal place of business at  
26 One Apple Park Way, Cupertino CA 95014.

1 **C. Agents and co-conspirators**

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

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

11  
12 **V. FACTUAL ALLEGATIONS**

13 19. Apple, currently the second largest company in the world by market capitalization, is an  
14 electronics and media company that designs, manufactures, and sells software and hardware  
15 technology products. Every second, Apple sells seven iPhones. Apple owns and operates the App  
16 Store, the primary platform for app developers to offer, and users to download, purchase, and subscribe  
17 to, mobile applications that run on Apple’s mobile operating system, iOS. The App Store earns billions  
18 of dollars of revenue for Apple each year.

19 20. In or around April 2024, Apple released a suite of open-source generative-AI models  
20 called OpenELM. Just a few months later, in or around June 2024, Apple announced the development  
21 of its commercial artificial-intelligence platform, called Apple Intelligence. Apple Intelligence includes  
22 multiple generative-AI models and related tools and technologies. The day after Apple officially  
23 introduced Apple Intelligence—what one reporter likened to ushering in a “new era”—the company  
24 gained more than \$200 billion in value: “the single most lucrative day in the history of the company.”

25 21. In January 2025, the company reported its “best quarter ever” with revenue of \$124.3  
26 billion, twice citing Apple Intelligence in its press release announcing the same and deeming it a part  
27 of the company’s “best-ever lineup of products and services.” The technology is integrated across  
28 Apple’s products, including iPhones, and is intended to “make[] apps and experiences even better and

1 more personal.” Analysts contend that Apple Intelligence could add \$4 trillion to the company’s market  
2 capitalization.

3 22. To train its generative-AI models, including the models that are part of Apple  
4 Intelligence, Apple first amassed an enormous library of data. Part of Apple’s data library includes  
5 copyrighted works—including books created by Plaintiffs—that were copied without author or owner  
6 consent, credit, or compensation. In essence, Apple has created its own internal shadow library for the  
7 purpose of training current and future LLMs that are and will be integrated into Apple software, Apple  
8 apps, and third-party apps that are offered on the Apple App Store and run on Apple hardware.

9 23. Apple has not attempted to pay these authors and owners for their contributions to this  
10 potentially lucrative venture. Apple did not seek licenses to copy and use the copyrighted books  
11 provided to its models. Instead, it intentionally evaded payment by using books already compiled in  
12 pirated datasets.

13 **A. How large language models work.**

14 24. Artificial intelligence denotes software that is designed to algorithmically create an  
15 illusion of human reasoning or inference, often using statistical and mathematical methods. While  
16 called “artificial intelligence,” these systems largely derive their value from human intelligence, as  
17 embodied within (among other things) the massive number of copyrighted works that are used to train  
18 them.

19 25. As noted above, though an LLM is a software program, it is not created the way most  
20 software programs are—by human software programmers writing code. Rather, LLMs like Apple’s are  
21 trained by copying an enormous quantity of textual works and then feeding these copies into the  
22 model. Typically, an LLM developer like Apple starts by identifying and copying its initial raw  
23 datasets containing such works. Then, these raw datasets are filtered and cleaned before they are used  
24 as model training data. For example, model developers often create software filters to remove certain  
25 categories of data that are undesirable for training, such as personally identifying information,  
26 copyright and license notices, or profanity and unsafe material. During this process, LLM developers  
27 often create derivative datasets that are subsets or altered versions of the original raw datasets.

28 Importantly, because these additional derivative datasets also contain copies of parts of or entire textual

1 works from the original raw datasets, including copyrighted works, the process of filtering and  
2 cleaning raw datasets to create usable training data for an LLM necessarily requires the creation of  
3 additional copies of the textual works in the raw datasets. Furthermore, because this gathering and  
4 processing of data is considered a strategic investment, the model developer often retains the datasets  
5 in a private data library, to remain available for future research and reuse.

6         26. The corpus of text ultimately used to train a model is called the “training dataset.”  
7 Datasets may have multiple uses during the development process of an LLM—for instance, checking  
8 the performance of a model at intermediate stages of training, in an iterative development process.  
9 Therefore, even if a dataset is not part of a model’s final training dataset, it may still be used in a  
10 model’s development process and retained in a data library for use in development of a future model.

11         27. During training, the LLM copies and ingests each textual work in the training dataset  
12 and extracts protected expression from it. The LLM progressively adjusts its output to more closely  
13 approximate the protected expression copied from the training dataset. The LLM records the results of  
14 this process in a large set of numbers called “weights” that are stored within the model. These weights  
15 are entirely and uniquely derived from the protected expression in the training dataset. Generally, the  
16 more data the LLM copies during training, the better the LLM’s ability to simulate the protected  
17 expression within that data as part of the LLM’s output. Moreover, generally speaking, the more  
18 diverse the training data is in prose, format, or subject matter, and the higher the quality of the written  
19 works that are within the training data, the better the LLM will be at producing high quality textual  
20 output. As a result, high-quality textual works, including longform creative works, are highly valuable  
21 to companies seeking to compile AI training datasets.

22         28. Once the LLM has copied and ingested the textual works in the training dataset and  
23 converted the protected expression into stored weights, the LLM can emit convincing simulations of  
24 natural written language in response to user prompts. Whenever an LLM generates a response to a user  
25 prompt, it is performing a computation that relies on these stored weights, with the goal of imitating  
26 the protected expression ingested from the training dataset.

1           29.     A significant portion of the material in the datasets that Apple downloaded and used for  
2 the purposes of training its LLMs consists of copyrighted works—including books written by Plaintiffs  
3 and Class Members—that Apple copied without consent and without providing credit or compensation.

4 **B.     The OpenELM language models were trained on copyrighted works.**

5           30.     In April 2024, Apple first announced the availability of the OpenELM language models  
6 on its website: “[W]e release OpenELM, a state-of-the-art open language model. OpenELM uses a  
7 layer-wise scaling strategy to efficiently allocate parameters within each layer of the transformer  
8 model, leading to enhanced accuracy.”

9           31.     The set of OpenELM language models released in April 2024 included variants called  
10 OpenELM-270M, OpenELM-450M, OpenELM-1\_1B, and OpenELM-3B. The main difference  
11 between these models is the parameter size; a larger parameter size means the model can store more  
12 weights and perform more complex tasks (requiring more computing power). For instance, Apple’s  
13 OpenELM-3B language model is so named because the model stores three billion (“3B”) weights  
14 derived from protected expression found in its training dataset.

15           32.     Each OpenELM model is hosted on a website called Hugging Face, where it has a  
16 “model card,” a file accompanying an AI model that typically describes the model, its intended uses  
17 and limitations, its training parameters, and the training dataset used to train the model. The model card  
18 for each OpenELM model states: “Our pre-training dataset contains RefinedWeb, deduplicated PILE, a  
19 subset of RedPajama, and a subset of Dolma v1.6, totaling approximately 1.8 trillion tokens.”

20           33.     Apple’s GitHub repository confirms that “OpenELM was pretrained on public datasets.  
21 Specifically, our pre-training dataset contains RefinedWeb, PILE, a subset of RedPajama, and a subset  
22 of Dolma v1.6.”

34. Apple also published a paper about OpenELM (“OpenELM Paper”). In a table called “Dataset used for pre-training OpenELM,” shown at right, Apple reveals that a large quantity of training data comes from the “Books” subset of a dataset called “RedPajama.” The OpenELM Paper does not further describe the contents of the RedPajama dataset.

Source	Subset	Tokens
RefinedWeb		665 B
	Github	59 B
	Books	26 B
	ArXiv	28 B
	Wikipedia	24 B
RedPajama	StackExchange	20 B
	C4	175 B
PILE		207 B
Dolma	The Stack	411 B
	Reddit	89 B
	PeS2o	70 B
	Project Gutenberg	6 B
	Wikipedia + Wikibooks	4.3 B

Table 2. Dataset used for pre-training OpenELM.

35. But information about the RedPajama dataset is available elsewhere. The RedPajama dataset is hosted on Hugging Face. According to the documentation for the RedPajama dataset that was available there until around April 2024, its “Books” component is a copy of the “Books3 dataset” that is “downloaded from Huggingface [sic]” when a user runs the script that automatically assembles the RedPajama dataset. Therefore, anyone who used the “Books” subset of the RedPajama dataset for training an AI model used a copy of the Books3 dataset. The documentation for the RedPajama dataset does not further describe the contents of Books3.

36. A paper describing the RedPajama dataset (the “RedPajama Paper”) further states that Apple’s OpenELM models were trained on RedPajama-V1, which was described as “a publicly available, fully open, best-effort reproduction of the training data . . . used to train the first iteration of LLaMA family of models.” The RedPajama Paper’s recitation of the components of the LLaMA training dataset includes six of the same components (with the same token counts) that are listed in Apple’s paper on OpenELM: GitHub, Books, ArXiv, Wikipedia, StackExchange, and C4. The RedPajama Paper later confirms that the “Books” component, with 26B tokens, included Books3.

37. A description of the contents of the Books3 datasets is available elsewhere. Books3 is a component of a separate AI training dataset called the Pile that was curated by a research organization called EleutherAI. In December 2020, EleutherAI introduced this dataset in a paper called “The Pile: An 800GB Dataset of Diverse Text for Language Modeling” (“The Pile Paper”). This paper describes the contents of Books3:

1 Books3 is a dataset of books derived from a copy of the contents of the  
2 Bibliotik private tracker ... Bibliotik consists of a mix of fiction and  
3 nonfiction books and is almost an order of magnitude larger than our  
4 next largest book dataset ... We included Bibliotik because books are  
invaluable for long-range context modeling research and coherent  
storytelling.

5 38. Similarly, Shawn Presser, who compiled the Books3 dataset, publicly confirmed in 2020  
6 that that dataset is composed of “all of Bibliotik” and includes .txt files containing the full text of  
7 approximately 196,640 books.

8 39. Bibliotik is one of several notorious “shadow library” websites. Others include Library  
9 Genesis (aka LibGen, Z-Library, or B-ok), Sci-Hub, and Anna’s Archive. The AI-training community  
10 has long been interested in these shadow libraries because they host and distribute vast quantities of  
11 unlicensed copyrighted material. For that reason, these shadow libraries violate the U.S. Copyright Act.

12 40. The 196,640 books in the Books3 dataset exist in .txt file format. A .txt file (pronounced  
13 a “text” file) is a simple file format that stores text data without any formatting, fonts, or images.  
14 Accordingly, the Books3 dataset consists of the text of the underlying 196,640 books.

15 41. By using the entire text of each book in Books3, Apple used copies of entire works to  
16 train its OpenELM models.

17 42. Plaintiffs’ Infringed Works are among the works in the Books3 dataset.

18 43. Until October 2023, the Books3 dataset was available from Hugging Face. At that time,  
19 the Books3 dataset was removed with a message that it “is defunct and no longer accessible due to  
20 reported copyright infringement.”

21 44. Presser himself has acknowledged that “we almost didn’t release the data sets at all  
22 because of copyright concerns.”

23 45. Before October 2023, anyone who used the “Books” subset of the RedPajama dataset  
24 for training necessarily made a copy of the Books3 dataset. Based on the information revealed in the  
25 OpenELM research paper and on its model card on Hugging Face, this includes Apple.

26 46. In sum, Apple has admitted training its OpenELM large language models on a copy of  
27 the “Books” subset of the RedPajama dataset, which in turn is a copy of the Books3 dataset. Therefore,  
28 Apple trained its OpenELM models on a copy of Books3, a known body of pirated books.

1 47. Because Plaintiffs' Infringed Works are part of Books3, Apple trained OpenELM on one  
2 or more copies of the Infringed Works and directly infringed Plaintiffs' copyrights along with the  
3 copyrights of the Class.

4 **C. The Apple Intelligence Foundation Language Models were trained on copyrighted works.**

5 48. In June 2024, Apple announced its commercial AI software platform, called Apple  
6 Intelligence. Apple Intelligence includes several AI models developed by Apple. Several of these  
7 models are called the *Apple Intelligence Foundation Language Models*.

8 49. The Foundation Language Models were described in a paper of the same name released  
9 by Apple on July 29, 2024 (the "FLM Paper"). According to that paper, the models are multimodal,  
10 meaning they can process, understand, and generate multiple types of data, including text and images.  
11 It follows that the Foundation Language Models were trained on image datasets as well as sets of text  
12 data.

13 50. The adjective *foundation* is commonly used to describe AI models that have broad  
14 capabilities to perform a wide variety of tasks. Consistent with this, Apple describes its Foundation  
15 Language Models as "highly capable in tasks like language understanding, instruction following,  
16 reasoning, writing, and tool use ... These foundation models are at the heart of Apple Intelligence."

17 51. The FLM Paper emphasizes the special importance of a foundation model's capacity to  
18 write: "[w]riting is one of the most critical abilities for large language models to have, as it empowers  
19 various downstream use[s]."

20 52. In the FLM Paper, Apple identifies two separate foundation language models:  
21 *AFM-server* and *AFM-on-device*. The *AFM-server* model is a larger model that is intended for use  
22 through an Apple-operated cloud service called Private Cloud Compute. The *AFM-on-device* model,  
23 by contrast, is intended to be small enough to be used directly on Apple devices (e.g., iPhones and  
24 laptops). According to the FLM Paper, the *AFM-on-device* model is "initialize[d] ... from a pruned  
25 6.4B model (trained from scratch **using the same recipe as AFM-server.**)" (emphasis added). This  
26 means that both the *AFM-server* and *AFM-on-device* models are trained on the same corpus of training  
27 data.

1           53. In the FLM Paper, Apple reveals three sources of training data: “data we have licensed  
2 from publishers, curated publicly available or open-sourced datasets, and publicly available  
3 information crawled by our web-crawler, Applebot.”

4           54. In describing the first source of training data—“data we have licensed”—Apple says  
5 only that it “identif[ied] and license[d] a **limited amount** of high-quality data from publishers”  
6 (emphasis added). In addition to being comparatively “limited” in quantity, Apple does not use this  
7 licensed data during the main phase of training the Foundation Language Models—what Apple calls  
8 “core pre-training”—but during a subsequent phase called “continued pre-training.”

9           55. As to its second source of training data—“publicly-available or open-sourced  
10 datasets”—Apple does not elaborate on the specific datasets used, saying only, “We evaluated and  
11 selected a number of high-quality publicly-available datasets with licenses that permit use for training  
12 language models.” Apple then filtered the datasets to remove personally identifiable information,  
13 copyright management information, and other undesired text before including them in the pre-training  
14 mixture.

15           56. In the parlance of AI training datasets, Apple’s phrase “publicly available” is one  
16 commonly used to falsely conjure up the idea of works made publicly available by the author. But the  
17 “public” nature of a dataset does not mean that the data collected in the dataset was obtained lawfully  
18 or that the party providing copies of the dataset has authority to extend a valid license to use the  
19 underlying copyrighted works. When companies describing AI training datasets use the phrase  
20 “publicly available” to describe those datasets, it has in practice meant only that the datasets were  
21 downloaded somewhere from the public internet, which contains a vast number of copyrighted works  
22 by authors and owners who have not granted a license for reproduction. There is a name for this kind  
23 of downloading of copies: copyright infringement. There is also a name for the infringing copies:  
24 pirated works.

25           57. For instance, Apple described its training data for OpenELM, including data from the  
26 Pile, as “public datasets,” even though the Pile contains an extraordinary number of pirated copies of  
27 copyrighted works. As another example, Meta Platforms also trained its LLaMA language models on  
28 Books3, which it described as a “publicly available dataset for training large language models” despite

1 the fact that none of the authors or owners whose works appear in Books3 ever consented to having  
2 their works included. Books3 was “publicly available” only in the limited sense that at one time, it  
3 could be acquired by anyone with an internet connection.

4 58. Similarly, in the context of AI training datasets, Apple’s phrase “open source” is  
5 commonly used to falsely conjure up the idea of works made available by the author under a  
6 permissive copyright license (e.g., a Creative Commons license). In practice, what it really means is  
7 copies of works made freely available by someone other than the author, without the author’s  
8 permission. Again, these are just pirated works.

9 59. For instance, EleutherAI, the group that created The Pile—the dataset that included  
10 Books3—described it as a “diverse, open source language modelling data set” even though the  
11 copyrighted works in the Books3 portion were included without authors’ and owners’ consent. Only  
12 the copyright holder can offer their copyrighted work to the public under an open-source license. A  
13 third party cannot usurp that right.

14 60. Therefore, when Apple says that a major source of the training data for its Foundation  
15 Language Models is “publicly-available or open-sourced datasets,” this should, on information and  
16 belief, be read to mean “certain pirated works or certain other pirated works.” Because Books3 has  
17 been described by people in the AI industry as a “publicly available” or “open-sourced” dataset, and  
18 because Apple already had a copy of Books3 that it had used for training its OpenELM models, on  
19 information and belief, Apple’s reference to “publicly-available or open-sourced datasets” includes  
20 Books3, and Apple therefore included Books3 in its private data library, which formed the basis for the  
21 training dataset for its Foundation Language Models. Thus, it is also, on information and belief, the  
22 case that the “curated publicly available or open-sourced datasets” that Apple copied to create a private  
23 data library for the purposes of training its Foundation Language Models (and future models) contain  
24 copyrighted material, including Plaintiffs’ and Class Members’ copyrighted works.

25 61. Plaintiffs’ Infringed Works are part of Books3. It follows that, on information and  
26 belief, Apple trained its Foundation Language Models on one or more copies of the Infringed Works,  
27 thereby directly infringing the copyrights of the Plaintiffs. On information and belief, Apple has  
28

1 created a permanent AI training-data library containing copies of all of its “curated,” purportedly  
2 “publicly-available or open-sourced” datasets, like Books3, in expectation of training future models.

3 62. As to its third source of training data—web pages crawled by Applebot—Apple says,  
4 “we crawl publicly available information using our web crawler, Applebot ... and respect the rights of  
5 web publishers to opt out of Applebot.” While Apple officially acknowledged Applebot’s crawling of  
6 webpages in mid-2015, its activity had been spotted in logs as early as November 2014. Around June  
7 2024, Apple revealed that it was using Applebot-scraped data for training its AI models. In response to  
8 this disclosure, by August 2024, numerous major commercial web publishers had chosen to opt out of  
9 Applebot training.

10 63. But Apple’s Foundation Language Models had necessarily been trained well before the  
11 release of the FLM Paper describing them in July 2024. For that reason, Apple’s disclosure in June  
12 2024 that it was using Applebot data to train language models came too late for any of these opt-outs to  
13 matter. Apple had already scraped the data and trained language models with it. On information and  
14 belief, Apple has retained copies of all Applebot data scraped before this wave of opt-outs, in  
15 expectation of training future models, as part of its AI training-data library.

16 64. In a November 2024 paper by George Wukoson and Joey Fortuna called “The  
17 Predominant Use of High-Authority Commercial Web Publisher Content to Train Leading LLMs,” the  
18 authors studied LLM training datasets made by algorithmically filtering scraped web pages. The  
19 authors concluded that such “datasets are disproportionately composed of high-quality content owned  
20 by commercial publishers of news and media websites.” In turn, this material is often covered by  
21 registered copyrights. Thus, on information and belief, the parts of Apple’s private shadow library and  
22 training datasets that come from filtered Applebot pages includes copyrighted works from commercial  
23 news and media websites.

24 65. The shadow libraries that host millions of unlicensed copyrighted books are also part of  
25 the “publicly available information” reachable by a web scraper like Applebot. Hence, on information  
26 and belief, part of Apple’s training-data library is sourced from shadow libraries via Applebot.  
27 On information and belief, Apple obscures the libraries and datasets used to train its Foundation  
28 Language Models to conceal its use of copyrighted materials. On information and belief, Apple’s

1 decision not to disclose the contents of its internal data libraries, or the training datasets for its  
2 Foundation Language Models, stems in part from the fact that Apple was the subject of negative press  
3 for using a subset of data from The Pile, containing captions from thousands of YouTube videos, at  
4 least in training the OpenELM models and, on information and belief, also in training the Foundation  
5 Language Models.

6 **D. Apple’s models for data classification were trained on copyrighted works.**

7 66. In the FLM Paper, Apple says that Applebot pages are “processed by a pipeline which  
8 performs quality filtering ... using heuristics and model-based classifiers.” In this context, the term  
9 “model-based classifier” refers to a separate AI model that has been trained to algorithmically rate the  
10 quality of and sort scraped web pages. On information and belief, the models that underlie these  
11 model-based classifiers were themselves trained on datasets that include unlicensed copyrighted works.

12 **E. All of Apple’s LLMs and model-based classifiers likely used multiple shadow libraries.**

13 67. As copyright-infringement litigation has proceeded against other companies training  
14 LLMs, it has become clear that these companies downloaded, torrented, and/or used more pirated  
15 works and shadow libraries than was publicly known at the time of the complaint. Across these suits,  
16 the additional shadow libraries implicated include Anna’s Archive, Pirate Library Mirror (“PiLiMi”),  
17 Z-Library (also known as B-ok), LibGen, and Sci-Hub. *See, e.g.*, Joint Letter Brief, *In re OpenAI*  
18 *ChatGPT Litig.*, No. 3:23-cv-03223-AMO (N.D. Cal. Jan 17, 2025) (Dkt. No. 254, at 1, 5)  
19 (acknowledging OpenAI’s production of datasets “LibGen 1” and “LibGen 2”); First Consolidated  
20 Amended Complaint, *Nazemian, et al. v. NVIDIA*, No. 4:24-cv-01454-JST (N.D. Cal. Jan. 16, 2026)  
21 (Dkt. No. 235, at 12-13) (describing communications between NVIDIA and Anna’s Archive, and  
22 alleging pirating from LibGen, Z-Library, Pirate Library Mirror, and Sci-Hub); Order Denying the  
23 Plaintiffs’ Motion for Partial Summary Judgment and Granting Meta’s Cross-Motion for Partial  
24 Summary Judgment, *Kadrey, et al. v. Meta*, No. 3:23-cv-03417-VC (N.D. Cal. June 25, 2025) (Dkt.  
25 No. 598, at 11) (“In early 2024, Meta also downloaded Anna’s Archive, a compilation of shadow  
26 libraries including LibGen, Z-Library, and others.”); Order on Fair Use, *Bartz, et al. v. Anthropic*, No.  
27 3:24-cv-05417-WHA (N.D. Cal. June 23, 2025) (Dkt. No. 231, at 3) (noting that Anthropic  
28 downloaded millions of copies of books from LibGen and PiLiMi). According to the administrators of

1 Anna’s Archive, “virtually all major companies building LLMs contacted us to train on our data. . . We  
2 have given high-speed access to about 30 companies.” The piracy is so widespread as to constitute a *de*  
3 *facto* industry practice.

4 68. Apple, like the companies listed above, is also a major company competing to build  
5 LLMs. As there was with the companies above, there is publicly available information indicating that  
6 Apple downloaded copyrighted works and used them in the multiple ways alleged herein. Thus, on  
7 information and belief, Apple, like the companies listed above, also very likely accessed one or more  
8 of these notorious shadow libraries to obtain large volumes of illegal copies of copyrighted books. Like  
9 other major technology companies that did the same thing, Apple, on information and belief, did so  
10 while avoiding paying the customary fee for books and without permission from the copyright holder.

11 **F. Apple integrates the Foundation Language Models, via Apple Intelligence, into Apple’s**  
12 **commercial products.**

13 69. Apple has integrated its Apple Intelligence products, including its Foundation Language  
14 Models, into a wide variety of commercially available products, including versions of the iPhone 15,  
15 16, and 17; the iPad, iPad mini, and iPad Pro; the MacBook Air, MacBook Pro, iMac, Mac mini, and  
16 Mac Studio; and Apple Vision Pro.

17 70. In addition, Apple has offered Apple Intelligence, and the FLM in particular, to  
18 developers of Apple device software applications, or “apps,” as a way to “create new intelligence  
19 features” using “the on-device large language model” at the core of Apple Intelligence. Apps are  
20 software applications, designed to run on devices such as smartphones, tablets, and other devices, that  
21 permit end-users to accomplish certain tasks. With many apps available on Apple devices, the end-  
22 users are typically consumers. According to Apple’s public-facing site for app developers, the  
23 “Foundation Models framework provides” developers “access to Apple’s on-device large language  
24 model that powers Apple Intelligence,” and can be used by developers to create apps that permit  
25 consumers to “generat[e] creative content” and complete “a diverse range of text generation tasks.”  
26 Apple’s site specifically highlights the ability of Apple’s FLM on-device model to “[c]ompose creative  
27 writing,” even including, as an example, the following prompt an app user could enter: “Generate a  
28

1 short bedtime story about a fox.” Additionally, Apple has released code and a machine-learning  
2 framework to assist development of iOS apps that integrate the OpenELM models.

3 71. Unsurprisingly, within several months of the Foundation Language Models becoming  
4 available to Apple developers for iOS 26, apps have appeared in the Apple App Store that allow users  
5 to generate output that could, in the short- or long-term, compete with copyrighted works. For  
6 example, Apple touted the app “Lil Artist” as one app relying on FLM, describing it as “combin[ing]  
7 the capabilities of the Foundation Models framework and the ImageCreator API to customize  
8 illustrated stories for children.” Another app available through the Apple App Store, called “Locally,”  
9 includes a “write professionally” button that encourages creative, long-form content generation. These  
10 apps, similar apps that currently exist, and similar apps that will inevitably be created in the near future  
11 using the FLM models or an OpenELM model, can be used to generate output that could compete  
12 with, dilute the market for, or simply be perceived and consumed by consumers as substitutes for  
13 Plaintiffs’ and Class Members’ copyrighted works. Apple earns revenue from the distribution and use  
14 of these apps via its App Store.

15 72. Apple has reportedly explored a paid tier for users of its Apple Intelligence products,  
16 with some early reports suggesting it may do so as early as 2027.

17 **G. Apple’s conduct impairs the market for Plaintiffs’ and Class Members’ works.**

18 73. Apple has neither paid nor sought permission from Plaintiffs for the use of their  
19 copyrighted works. Instead of doing so, Apple downloaded, torrented, scraped, or otherwise copied  
20 vast quantities of copyrighted works—including illegally compiled datasets such as Books3—that  
21 included Plaintiffs’ works like the Infringed Works.

22 74. In so doing, Apple has, first, deprived Plaintiffs of the revenue that would have been  
23 generated had Apple approached Plaintiffs or their licensing agents directly to license copies of their  
24 works. Plaintiffs and similarly situated creators previously licensed their work for their own  
25 commercial uses.

26 75. There are numerous examples of publicly reported AI licensing deals. For example,  
27 Microsoft reportedly reached an agreement with HarperCollins to license a select number of non-  
28 fiction titles to train AI models for just three years; compensation for participating authors, who have

1 to opt-in before their works can be licensed, is \$2,500 per work. Wiley likewise reported in 2025 that it  
2 had executed multiple AI licensing agreements totaling \$40 million. Myriad licensing systems have  
3 been launched and are continuing to develop, including the Copyright Clearance Center’s collective AI  
4 licensing scheme and the Created by Humans licensing platform. Further, several AI dataset licensing  
5 companies have formed a trade group called the Dataset Providers Alliance in order to promote  
6 standardization and transparency in the licensing of intellectual property content for AI and machine  
7 learning datasets. Currently, some researchers estimate, the AI training license market is valued at  
8 approximately \$2.5 billion; within a decade, it may close in on nearly \$30 billion.

9 76. Apple itself understands the value of copyrighted works and the market that exists for  
10 paying creators to use their works for training. For instance, it struck an agreement with Shutterstock to  
11 “use hundreds of millions of images, videos and music files” valued between an estimated \$25 to \$50  
12 million. Similarly, Apple has contacted news organizations like Condé Nast, NBC News, and IAC to  
13 license news article archives. Nonetheless, Apple has not compensated Plaintiffs’ and Class Members  
14 whose works it copied and used in training its models.

15 77. Apple has not only directly deprived Plaintiffs and Class Members of the licensing  
16 revenues Plaintiffs and Class Members could have earned from Apple for Apple’s use of their works to  
17 train their LLMs. Apple has also contributed, through its conduct, to the impairment of the emergence  
18 of lawful licensing regimes. Such nascent licensing regimes are more successful if large corporate  
19 licensees participate early and at scale, contributing significant licensing revenues and stabilizing the  
20 demand side of such a licensing market. Apple’s choice to instead copy pirated works on a massive  
21 scale has, merely through the absence of a corporate cornerstone of such a market, blunted the  
22 development of an AI training dataset licensing regime.

23 78. Apple’s conduct has also depressed the demand for licensing from other potential  
24 licensees by normalizing and encouraging other potential licensees not to license works for AI training.  
25 Apple has significant consumer brand loyalty and is one of the most valuable companies in the world.  
26 Smaller potential licensees are less likely to pursue licensing if they see that Apple, with its cadre of  
27 retained and in-house lawyers, has determined it does not need to do so. Smaller potential licensees are  
28 also less likely to pursue licensing if they believe they cannot do so and remain competitive with larger

1 companies with more capital (like Apple) who have declined to pay licensing fees and instead flouted  
2 copyright law. Finally, smaller potential licensees are less likely to pursue licensing if Apple’s conduct  
3 has normalized the use of pirated works, which use Apple has largely attempted to hide from  
4 consumers.

5 79. Second, Apple’s unauthorized use of Plaintiffs’ copyrighted works to train the  
6 OpenELM models and Foundation Language Models has also caused and threatens to cause substantial  
7 harm to the actual and potential markets for those works. As described above, Apple has developed,  
8 designed, created, maintained, and/or distributed functionalities and applications that are designed to,  
9 and can, generate works that imitate or otherwise resemble the Infringed Works. Apple has also  
10 encouraged and facilitated third parties’ design, creation, maintenance, and/or distribution of relevantly  
11 similar applications. All of these unauthorized uses of the Infringed Works creates a risk of market  
12 dilution, which may “lead to a loss of sales” by “harm[ing] the market for access to those works.”  
13 Works generated using the OpenELM models or the Foundation Language Models will inevitably start  
14 competing with Plaintiffs’ works (like the Infringed Works) and ultimately dilute royalty pools as AI-  
15 generated output increasingly floods the market. Already, “low-quality sham ‘books’” created by other  
16 LLMs have begun overwhelming the market as scammers generate “unauthorized ‘biographies’ of  
17 authors that are simply AI-generated rehashings of their lives, often based on autobiographical works.”  
18 Other scams include “companion books” that summarize the key points from the original novel, with  
19 “little to no original analysis or commentary and are meant only to confuse consumers and skim sales  
20 off of the real books.” These works have already entered book marketplaces like Amazon.

21 80. Apple’s models—whether integrated in Apple Intelligence, an Apple-developed app, or  
22 a third-party-developed app—are designed to, and can, generate outputs that may substitute the kinds  
23 of expressive written work that Plaintiffs are hired to produce, potentially diminishing demand for  
24 books and human-produced stories. Plaintiffs face potential ongoing harms through lost publication  
25 opportunities, reduced recognition, and lost sales, among other harms.

26 81. This type of substitution of human-created works for AI-generated works has already  
27 been documented in creative markets. For example, the *New York Times* reported in a recent front-page  
28 story: “Publishers and authors worry that books by real writers are getting lost in the sea of digital slop,

1 as A.I.-enabled novels flood the market.” According to one best-selling novelist quoted in that article:  
2 “It bogs down the publishing ecosystem that we all rely on to make a living . . . It makes it difficult for  
3 newer authors to be discovered, because the swamp is teeming [sic] with crap.” Although some stigma  
4 currently attaches to known AI-generated novels, another writer predicts, “[e]ventually . . . readers will  
5 not care.” As an additional example, researchers have recently shown that some readers prefer the  
6 outputs of AI models trained on copyrighted works over the work of human writers.

7 82. Additionally, in a 2025 working paper analyzing the impact of generative AI on the  
8 stock image market, the authors describe their findings as follows:

9 In this paper, we ask three questions with the aims of shedding light on how GenAI will  
10 change markets for creative goods and informing the active debate on copyright and fair  
11 use.

12 **First, to what extent do GenAI goods compete with—or become substitutes for—**  
13 **non-GenAI goods? We provide causal evidence that consumers view GenAI-**  
14 **produced goods as substitutes for traditional creative goods. Thus GenAI goods**  
15 **are competitive with non-GenAI goods.** Further, the market expands as GenAI goods  
16 enter the market, with GenAI goods consumption more than replacing non-GenAI  
17 goods.

18 **Next, does GenAI result in crowd-out of non-GenAI firms and goods?** Policymakers  
19 in particular have expressed this concern given the substantial differences in production  
20 costs brought by GenAI. **Our results show that GenAI does crowd out non-GenAI**  
21 **firms and goods. Across almost all markets, we see substantial entry of GenAI**  
22 **firms and exit of non-GenAI firms, with net production expansion, but a decrease**  
23 **in any given product’s probability of sale. We document that this concern is**  
24 **particularly relevant for niche markets, where we find production tipping almost**  
25 **completely to GenAI. . . .**

26 . . .  
27 The substitution to non-GenAI content, and significant exit of non-GenAI artists,  
28 provide evidence of crowd-out of non-GenAI content. **In the long run, this may lead**  
29 **to a decline in the production of novel original content for model training, and**  
30 **even to a decline in non-GenAI content altogether. Further, we document a**  
31 **reduction in sales rates for artists.** While these reductions may be partially  
32 compensated for by the increase in production, not all artists, and in particular not all  
33 non-GenAI artists, will be able to fully compensate via scale, and thus their long-run  
34 profitability may be at risk. **Thus, our results provide some validity to concerns that**  
35 **GenAI content will both replace non-GenAI content and greatly diminish the value**  
36 **of the original content, affecting the fourth pillar of fair use.**

1 **VI. CLASS ACTION ALLEGATIONS**

2 83. The “Class Period” as defined in this Complaint begins at least three years before  
3 September 5, 2025 and runs through the present.

4 84. As used here, the term “Apple Uses” refers collectively to Apple’s numerous separate  
5 infringing uses of Plaintiffs’ and Class Members’ works, including the Infringed Works, as set forth  
6 below:

7 A. *Acquiring Data via Pirated Sources:* Apple’s unauthorized reproduction of  
8 Plaintiffs’ and Class Members’ works, including the Infringed Works, as part of Apple’s  
9 downloading, torrenting, or other mode of initial acquisition of one or more datasets of pirated  
10 works, including Books3, Library Genesis, Z-Library, Anna’s Archive, Sci-Hub, Internet  
11 Archive, and/or Pirate Library Mirror;

12 B. *Acquiring Data via Applebot:* Apple’s unauthorized copying of Plaintiffs’ and  
13 Class Members’ works, including the Infringed Works, as part of Apple’s scraping of the public  
14 internet using Applebot or other methods;

15 C. *Curating:* Apple’s unauthorized reproduction and use of Plaintiffs’ and Class  
16 Members’ works, including the Infringed Works, when filtering, cleaning, and processing at  
17 least one pirated dataset, such as Books3, and components of such datasets, to create a curated,  
18 private shadow library of datasets for use in various stages of training LLMs;

19 D. *Training OpenELM Models:* Apple’s unauthorized reproduction and use of  
20 Plaintiffs’ and Class Members’ works, including the Infringed Works, when training its  
21 OpenELM models;

22 E. *Training FLMs:* Apple’s unauthorized reproduction and use of Plaintiffs’ and  
23 Class Members’ works, including the Infringed Works, when training its Foundation Language  
24 Models;

25 F. *Training Classifier Models:* Apple’s unauthorized reproduction and use of  
26 unlicensed copyrighted works when training classifier models; and

27 G. *Retaining for Private Data Library:* Apple’s unauthorized retention and use of  
28 Plaintiffs’ and Class Members’ works, including the Infringed Works, along with all the data

1 Apple has illegally downloaded and curated thus far, in the form of a private data library for use  
2 in models currently under development as well as future models—an AI data library that  
3 includes the Books3 dataset, which, in turn, includes the Infringed Works.

4 85. **Class definition.** Plaintiffs bring this action for damages and injunctive relief as a class  
5 action under Federal Rules of Civil Procedure 23(a), 23(b)(2), and 23(b)(3), on behalf of the following  
6 Class:

7 **All beneficial or legal owners of a registered United States copyright**  
8 **in any work that was 1) used in one or more of the Apple Uses**  
9 **during the Class Period, and 2) was registered with the United States**  
10 **Copyright Office a) within five years of the work’s publication and**  
11 **before its use in one or more of the Apple Uses, or b) within three**  
12 **months of publication.**

13 86. This Class definition excludes: (a) Defendant; (b) Any of Defendant’s parent  
14 companies, subsidiaries, and affiliates; (c) Any of Defendant’s officers, directors, management,  
15 employees, subsidiaries, affiliates, or agents; (d) All governmental entities; and (e) The judges and  
16 chambers staff in this case, as well as any members of their immediate families.

17 87. **Numerosity: Federal Rule of Civil Procedure 23(a)(1).** The Class Members are so  
18 numerous and geographically dispersed that individual joinder of all Class Members is impracticable.  
19 The exact number of Class Members is currently unknown to Plaintiffs, as this information is in  
20 Defendant’s exclusive control. On information and belief, there are more than several thousand  
21 members in the Class across the United States. Accordingly, joinder of all Class Members in  
22 prosecuting this action is impracticable.

23 88. The Class can be identified, in part, through tools that allow a user to search for web  
24 domains included in the RedPajama dataset, the Pile dataset, and other datasets used for one or more of  
25 the Apple Uses.

26 89. The Class can further be identified by obtaining business records maintained by Apple,  
27 including the content of its private pirated libraries and datasets used for both its OpenELM models  
28 and Foundation Language Models.

1           90.     **Typicality: Federal Rule of Civil Procedure 23(a)(3).** Plaintiffs’ claims are typical of  
2 the claims of Class Members because Plaintiffs and all members of the Class were damaged by the  
3 same course of conduct of Defendant. Further, the relief sought is common to all Class Members.

4           91.     **Adequacy of Representation: Federal Rule of Civil Procedure 23(a)(4).** Plaintiffs  
5 will fairly and adequately represent the interests of the members of the Class because Plaintiffs have  
6 experienced the same harms as the members of the Class and have no conflicts with any other  
7 members of the Class. Further, Plaintiffs have retained competent counsel who are experienced in  
8 litigating federal class actions and other complex litigation involving sophisticated, state-of-the art  
9 technology.

10          92.     **Commonality and Predominance: Federal Rules of Civil Procedure 23(a)(2) and**  
11 **23(b)(3).** Numerous questions of law and fact are common to each Class Member arising from  
12 Defendant’s conduct, including:

13           A.     Whether Apple acquired or downloaded reproductions of datasets, including but  
14 not limited to the RedPajama and Books3 datasets, that included Plaintiffs’ and Class Members’  
15 copyrighted works;

16           B.     Whether Apple copied Plaintiffs’ and Class Members’ works as part of Apple’s  
17 scraping of the public internet using Applebot or other methods;

18           C.     Whether Apple included Plaintiffs’ and Class Members’ works in any of the  
19 datasets used by Apple to train its OpenELM models and Foundation Language Models;

20           D.     Whether Apple’s inclusion of Plaintiffs’ and Class Members’ works in those  
21 datasets constituted or required the works’ reproduction by Apple;

22           E.     Whether Apple lacked authorization to reproduce copies of Plaintiffs’ and Class  
23 Members’ works;

24           F.     Whether Apple violated the copyrights of Plaintiffs and the Class when it  
25 acquired or downloaded copies of Plaintiffs’ and Class Members’ works and used them in  
26 training its OpenELM models;

1 G. Whether Apple violated the copyrights of Plaintiffs and the Class when it  
2 acquired or downloaded copies of Plaintiffs' and Class Members' works and used them in  
3 training its Foundation Language Models;

4 H. Whether Apple violated the copyrights of Plaintiffs and the Class when it  
5 acquired or downloaded copies of Plaintiffs' and Class Members' works and retained them in a  
6 curated, private shadow library for Apple's use in training other models;

7 I. Whether this Court should enjoin Defendant from engaging in the unlawful  
8 conduct alleged herein;

9 J. Whether any affirmative defense excuses Defendant's conduct, including the fair  
10 use doctrine; and

11 K. Whether Apple's infringement was willful.

12 93. These and other questions of law and fact are common to the Class and predominate  
13 over questions affecting Class Members on an individual basis.

14 94. **Predominance & Superiority: Federal Rule of Civil Procedure 23(b)(3).** Defendant  
15 has acted on grounds generally applicable to the Class. A class action is superior to alternatives for the  
16 fair and efficient resolution of this controversy. Allowing the claims to proceed on a class basis will  
17 eliminate the possibility of repetitive litigation. Further, injunctive relief is appropriate with respect to  
18 the entire Class.

19 95. **Risk of Prosecuting Separate Actions.** The alternative of separate actions by  
20 individual Class Members risks inconsistent adjudications and is an inefficient use of limited judicial  
21 resources.

22  
23 **VII. CLAIM**

24 **DIRECT COPYRIGHT INFRINGEMENT — 17 U.S.C. § 501**

25 96. Plaintiffs incorporate by reference all other allegations in this complaint.

26 97. As the owners of the registered copyrights in the Infringed Works and other copyrighted  
27 works, Plaintiffs and Class Members hold the exclusive rights to those works under 17 U.S.C. § 106.

28 Plaintiffs and the Class Members never authorized Apple to make copies of their Infringed Works and

1 other copyrighted works, make derivative works, publicly display copies (or derivative works), or  
2 distribute copies (or derivative works), or exploit any other right exclusively reserved to Plaintiffs and  
3 the Class Members under the U.S. Copyright Act.

4 98. In performing the Apple Uses, including by downloading the works and making  
5 additional copies during various stages of training and development of the OpenELM models and the  
6 Foundation Language Models, Apple made all its copies of the Infringed Works and other copyrighted  
7 works without Plaintiffs' or Class Members' permission, violating their exclusive rights under the U.S.  
8 Copyright Act. Indeed, "the person who copies the textbook from a pirate site has infringed already,  
9 full stop." *Bartz et al. v. Anthropic*, 787 F. Supp. 3d 1007, 1025 (N.D. Cal. 2025). Regardless of how  
10 Apple uses the works in its private training-data library in the future, this cannot negate that the initial  
11 copying of works sourced from shadow libraries infringed on Plaintiffs' and Class Members' exclusive  
12 rights.

13 99. Plaintiffs and Class Members have been injured by Apple's acts of direct copyright  
14 infringement of the Infringed Works. Plaintiffs and Class Members are entitled to statutory damages,  
15 actual damages, restitution of profits, destruction of the infringing copies and models, and other  
16 remedies provided by law.

17 100. Apple's violation of Plaintiffs' and Class Members' exclusive rights was willful. Apple  
18 knew the datasets it downloaded, copied, stored, and trained its LLMs on contained copyrighted works.

## 20 VIII. PRAYER FOR RELIEF

21 101. Plaintiffs demand judgment on their behalf and on behalf of the Class against each  
22 Defendant as follows:

23 A. Allowing this action to proceed as a class action, with Plaintiffs serving as Class  
24 Representatives, and with Plaintiffs' counsel as Class Counsel;

25 B. Awarding Plaintiffs and the Class statutory damages, compensatory damages,  
26 restitution, disgorgement, and any other relief that may be permitted by law or equity;

27 C. Permanently enjoining Defendant from the unlawful, unfair, and infringing  
28 conduct alleged herein;

1 D. Ordering destruction under 17 U.S.C. § 503(b) of all Apple LLMs that have  
2 ingested at least one copy of any of Plaintiffs’ and Class Members’ works;

3 E. Ordering destruction under 17 U.S.C. § 503(b) of all copies of those works that  
4 are maintained in Apple’s private libraries and datasets;

5 F. An award of costs, expenses, and attorneys’ fees as permitted by law; and

6 G. Such other or further relief as the Court may deem appropriate, just, and  
7 equitable.

8  
9 **IX. DEMAND FOR JURY TRIAL**

10 Plaintiffs demand a jury trial for all claims.

11  
12 DATED this 13th day of February, 2026.

13 KELLER ROHRBACK L.L.P.

14 By s/ William K. Dreher

15 William K. Dreher, *Pro Hac Vice*  
16 Benjamin Gould (SBN 250630)  
17 Derek W. Loeser, *Pro Hac Vice*  
18 Cari C. Laufenberg, *Pro Hac Vice*  
19 Chris N. Ryder, *Pro Hac Vice*  
20 Elizabeth W. Tarbell, *Pro Hac Vice*  
21 Samuel Rubinstein, *Pro Hac Vice*  
1201 Third Avenue, Suite 3400  
22 Seattle, WA 98101-3268  
23 (206) 623-1900  
24 Fax (206) 623-3384  
25 bgould@kellerrohrback.com  
wdreher@kellerrohrback.com  
26 dloeser@kellerrohrback.com  
cryder@kellerrohrback.com  
27 claufenberg@kellerrohrback.com  
28 etarbell@kellerrohrback.com  
srubinstein@kellerrohrback.com

SUSMAN GODFREY L.L.P.

By s/ Rohit D. Nath

Rohit D. Nath (SBN 316062)  
Justin A. Nelson, *Pro Hac Vice*

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

Alejandra C. Salinas, *Pro Hac Vice*  
Jordan Connors, *Pro Hac Vice*  
Eleanor Runde, *Pro Hac Vice*  
1900 Avenue of the Stars, Suite 1400  
Los Angeles, CA 90067  
(310) 789-3100  
rnath@susmangodfrey.com  
jnelson@susmangodfrey.com  
asalinas@susmangodfrey.com  
jconnors@susmangodfrey.com  
erunde@susmangodfrey.com

BUTTERICK LAW PC

By s/ Matthew Butterick  
Matthew Butterick (SBN 250953)  
1920 Hillhurst Avenue, #406  
Los Angeles, CA 90027  
mb@buttericklaw.com

Attorneys for Plaintiffs