

**LATHAM & WATKINS LLP**

Joseph R. Wetzel (SBN 238008)

*joe.wetzel@lw.com*

Andrew M. Gass (SBN 259694)

*andrew.gass@lw.com*

Brittany N. Lovejoy (SBN 286813)

*britt.lovejoy@lw.com*

505 Montgomery Street, Suite 2000

San Francisco, California 94111

Telephone: +1.415.391.0600

Sarang V. Damle (*pro hac vice*)

*sy.damle@lw.com*

555 Eleventh Street NW, Suite 1000

Washington, D.C. 20004

Telephone: +1.202.637.2200

Allison L. Stillman (*pro hac vice*)

*alli.stillman@lw.com*

1271 Avenue of the Americas

New York, New York 10020

Telephone: +1.212.906.1747

*Attorneys for Defendant Anthropic PBC*

**UNITED STATES DISTRICT COURT  
NORTHERN DISTRICT OF CALIFORNIA  
SAN JOSE DIVISION**

CONCORD MUSIC GROUP, INC., ET AL.,

Plaintiffs,

vs.

ANTHROPIC PBC,

Defendant.

Case No. 5:24-cv-03811-EKL

**DECLARATION OF JARED KAPLAN  
IN SUPPORT OF DEFENDANT'S  
OPPOSITION TO PLAINTIFFS' RENEWED  
MOTION FOR PRELIMINARY  
INJUNCTION**

Hon. Eumi K. Lee

**REDACTED VERSION**

***[Refiled Pursuant to Order at Dkt. 282; Originally filed on August 22, 2024 at Dkt. 209]***

1 I, Jared Kaplan, declare the following:

2 1. I am a co-founder and the Chief Science Officer of Defendant Anthropic PBC and  
3 submit this declaration in support of Defendant's Opposition to Plaintiffs' Renewed Motion for  
4 Preliminary Injunction. Unless stated otherwise, all facts stated herein are within my personal  
5 knowledge.

6 2. I have a B.S. in physics and mathematics from Stanford University and a Ph.D. in  
7 physics from Harvard University. I began my career as a theoretical physicist, with research  
8 interests in effective field theory, particle physics, cosmology, scattering amplitudes, and the  
9 conformal field theory bootstrap. I have been a professor in the Department of Physics and  
10 Astronomy at Johns Hopkins University since 2012. My work in theoretical physics led me to  
11 collaborate with computer scientists to conduct academic research into machine learning and the  
12 development of large language models. I have taught courses and published over 60 scholarly  
13 articles on those subjects.

14 3. I have been the Chief Science Officer of Anthropic since 2023, and I was  
15 previously a manager at Anthropic since its founding in 2021. I have extensive personal  
16 knowledge regarding Anthropic's business operations and practices, specifically: knowledge of  
17 how Anthropic's generative artificial intelligence model was created and trained, its technical  
18 capacities and processes, Anthropic's goals for the model, and Anthropic's business operations  
19 related to the model. If called upon, I would be willing to testify as set forth in this declaration.

20 4. I have read the Complaint, the original and renewed motions for preliminary  
21 injunction, and the declarations of Michael Smith, Benjamin Zhao, and Edward Newton-Rex  
22 filed by Plaintiffs in this action. In filing their Renewed Motion, I understand Plaintiffs have  
23 asked this Court to do two things: (1) order Anthropic to maintain its current guardrails that  
24 prevent its AI models from generating outputs that disseminate Plaintiffs' song lyrics; and (2)  
25 prohibit Anthropic from using the lyrics to songs allegedly owned or controlled by Plaintiffs to  
26 train any of Anthropic's future AI models.

**I. WHAT IS GENERATIVE ARTIFICIAL INTELLIGENCE?**

**A. Generative Artificial Intelligence**

5. Generative AI models are software designed to generate original text, images, and computer code at a level that approximates human capabilities.

6. One category of generative AI models is large language models (“LLMs”). Generally speaking, LLMs are text-based generative AI models that use deep learning techniques and train on extremely large data sets to understand, summarize, generate, and predict new content. In addition to generating new texts, LLMs can also translate texts, summarize or rewrite content, or engage in conversation with users. The purpose of generative AI models is to understand user queries and respond intelligently in a manner approximating a human assistant.

**B. Anthropic and the Claude Model**

7. Anthropic is an artificial intelligence safety and research company based in San Francisco, California, working to develop generative AI models.

8. Anthropic’s signature product is Claude, a series of general purpose LLMs. Although Anthropic was originally founded as a research company, Anthropic launched its first commercial model via an application programming interface (“API”) to certain business users in late 2022 following two years of research, development, and training. Anthropic released a public conversational interface, Claude.ai, as a consumer product in July 2023. An updated version of the model, Claude 2, was also released in July 2023. [REDACTED]

[REDACTED] The Claude 3 model family (Haiku, Sonnet, and Opus) was released on March 4, 2024, with Claude 3.5 Sonnet following on June 20, 2024. [REDACTED]

9. Like many LLMs, Claude interprets and responds to a vast variety of user queries, known as “prompts,” in an intelligent, human-like way. As a result of Anthropic’s tremendous research, development, and training of the model, Claude’s novel outputs are produced

1 algorithmically and occur automatically without real-time, human input. Claude is designed to  
2 use the learning from its training to facilitate the creation of new, original content.

3 10. Claude is a versatile LLM that can be used in many different contexts depending  
4 on a user's needs and provides novel, tailored outputs for users on myriad topics. Claude allows  
5 users to engage in general, open-ended, and back-and-forth conversations involving complex  
6 reasoning, creativity, thoughtful dialogue, software coding, or detailed content creation.

7 11. Depending on the nature of the user prompt, Claude can:

- 8 • Handle basic instructions and logical scenarios, including formatting outputs,  
9 following if-then statements, and making a series of logical evaluations in a  
10 single prompt.
- 11 • Take direction on tone and "personality" when providing outputs.
- 12 • Speak a variety of common human languages.
- 13 • Provide outputs in a variety of programming languages, including generating  
14 code based on natural language descriptions, troubleshooting and adding  
15 functionality to existing codebases, and translating code from one  
16 programming language to another.
- 17 • Process huge amounts of text from documents, emails, FAQs, chat transcripts,  
18 records, charts, or graphs and edit, rewrite, summarize, classify, extract  
19 structured data, draft Q&As, and complete other tasks based on the content.
- 20 • Interact with external tools and APIs to perform tasks, manipulate data, and  
21 provide dynamic and accurate responses.
- 22 • Provide helpful answers to user queries, with detailed background on  
23 technical, scientific, and cultural knowledge.

24 12. On average, hundreds of thousands of users engage with Claude.ai on a daily  
25 basis, which amounts to approximately millions of conversations per day. In practice, individual  
26 users typically rely on Claude to assist with original writing projects including editing, rewriting,  
27 summarizing, brainstorming ideas, and problem solving. Claude can assist users in  
28

1 straightforward tasks like drafting professional emails and job descriptions,<sup>1</sup> or in developing  
2 more complex, creative outputs, such as novels and poetry.<sup>2</sup>

3 13. In addition to individual use cases, businesses utilize Claude to support their  
4 operations and day-to-day needs, such as developing online tutoring programs, analyzing  
5 contracts, facilitating productivity and product management improvements, among other things.  
6 For example, one of Anthropic's partners, a leading provider of online education solutions,  
7 powers its online platform with Claude to help students achieve academic success by providing a  
8 conversational experience that matches the level of a human tutor. Legal infrastructure  
9 businesses use Claude to provide their customers with an AI platform capable of, for example,  
10 summarizing, analyzing, evaluating, and even drafting legal contracts. Asset management firms  
11 utilize Claude's capabilities to generate elaborate charts, compute financial indicators, and create  
12 summaries of the results.

13 14. A demonstrative video showing a business use case of Claude can be seen [here](#).  
14 Claude is also used as a springboard for creativity, such as acting as a writing partner and  
15 fleshing out plot points for an original story as seen [here](#), or interpreting and implementing ideas  
16 into code as seen [here](#).

17 15. Anthropic has not developed Claude, nor does it intend to develop Claude, to  
18 function simply as a resource for users to obtain copies of already-existing materials in response  
19 to queries. Rather, our goal is to generate *original* outputs that meet the full range of customer  
20 needs. As explained in more detail below, if Claude detects that a user is attempting to elicit an  
21 output that does not comply with Anthropic's Constitution (described below) and/or Anthropic's  
22 Terms of Use and Acceptable Use Policy (including the prohibition against generating infringing  
23 output), Claude will not provide the user with the requested output.

24 \_\_\_\_\_  
25 <sup>1</sup> *Introducing Claude*, Anthropic (Mar. 14, 2023), <https://www.anthropic.com/news/introducing-claude>.

26 <sup>2</sup> Anthropic, *Claude 3.5 Sonnet as a writing partner*, YOUTUBE (June 20, 2024),  
27 <https://www.youtube.com/watch?v=-dWfl7Dhb0o>; see also Anthropic, *Claude 3.5 Sonnet for*  
28 *sparkling creativity*, YOUTUBE (June 20, 2024),  
<https://www.youtube.com/watch?v=rHqk0ZGb6qo>.

1           16.     Such misuses of Claude are not the typical use case. After this lawsuit was filed,  
2 Anthropic collected data from users' accounts believed to be associated with Plaintiffs and their  
3 agents. Upon comparison of Plaintiffs' use of Claude versus an average user's use of Claude, it  
4 appears that Plaintiffs were in fact using Claude to continually prompt for lyrics of the songs in  
5 Exhibit A of their Complaint in a manner that does not remotely reflect typical (or permitted) use  
6 cases.

## 7 **II.     HOW CLAUDE IS DEVELOPED AND TRAINED**

### 8 **A.     Building and Training LLMs In General**

9           17.     Based on my experience in the field and at Anthropic, I am personally familiar  
10 with the standard process of building LLMs like Claude. To build an LLM like Claude,  
11 engineers first develop a "neural network," which is a computer program capable of studying  
12 enormous datasets of texts and extracting statistical and linguistic information about those texts,  
13 in an effort to gain a generalized understanding of human language. This neural network is the  
14 foundation for training any LLM and allows engineers to input the tremendous amount of  
15 materials needed to competently train LLMs.

16           18.     The neural network is shown a vast dataset consisting of hundreds of millions, if  
17 not billions, of pieces of content. This set of materials is known collectively as the "training  
18 corpus" or simply the "corpus."

19           19.     As an initial step, the corpus is preprocessed and broken down into "tokens,"  
20 which are basic semantic units that can be represented mathematically. Tokens can correspond  
21 to words, subwords, or characters. Once the materials in the training corpus are transformed  
22 from human-readable text into the tokens, the tokens are typically grouped into sequences (often  
23 around [REDACTED] tokens per group, although this can vary). These sequences are then shuffled to  
24 make random batches for training, which is meant to prevent the model from memorizing the  
25 corpus' original arrangements, or learning false relationships that only exist because of how the  
26 original corpus was assembled.

27           20.     The neural network then analyzes each token in the context of surrounding tokens  
28

1 to learn language patterns and the relationship between the words and phrases in the corpus.  
2 Through iterative training on large amounts of data, the neural network discovers patterns among  
3 words and phrases, extracting statistical information to develop a probabilistic but  
4 comprehensive map of how language works, without explicit human-defined rules. The model  
5 stores this statistical information in an organized set of numerical values called parameters,  
6 which include weights and biases (representing the connective strengths among neurons).  
7 During inference—the process of running new inputs through the trained model—the model uses  
8 these parameters to analyze patterns among tokens and generate text.

9       21. The goal of training an LLM is to develop a generalized understanding of  
10 language patterns—not to memorize specific content from the training corpus. This enables the  
11 LLM to understand a virtually limitless array of potential user prompts, and produce a broad  
12 range of novel outputs in response.

13       22. Trillions of tokens are required to properly train an LLM like Claude and to  
14 ultimately allow the LLM to recognize patterns in language to any useful degree. The volume is  
15 necessary to allow the LLM to gain a generalized understanding of how humans read and write  
16 text, so that it can interpret novel prompts and produce novel outputs.

17       23. In addition to a broad volume of data, any functional general purpose LLM  
18 requires a very diverse set of training data to gain a sufficiently precise understanding of the full  
19 range of human language and human experience. LLMs like Claude are general-purpose tools  
20 that need to be familiar with everything from how emails are drafted, how to write programs in  
21 C++, how many states are in the Union, and on and on. It follows that when texts are similar  
22 they are considered fungible for purposes of the model. For example, it does not ultimately  
23 matter what specific articles are included in the training data, so long as articles in general are  
24 included to diversify the training corpus.

25       24. In short, the more data tokens LLMs are able to process, the better they are able to  
26 represent a wide array of subject areas, generalize beyond the corpus, and improve the accuracy  
27 of their outputs.  
28

25. As a general matter, LLMs that are trained on fewer data tokens—e.g., an LLM trained solely on licensed or public domain data—will have lower performance and potentially fewer capabilities than models trained on trillions of tokens, like Claude. LLMs that are not trained on the same diversity of data, for example, may be limited in the types of tasks that they can perform, the same way that a human who has never learned how to write computer code—or has only taken Computer Science 101—may not know how to build or debug a complex computer program.

26. While LLMs trained on smaller corpora of data can be effective at a narrow set of tasks for which they have been specifically trained, they are not comparable to large, general-purpose LLMs like Claude. For instance, KL3M, the LLM mentioned by Plaintiffs’ expert Edward Newton-Rex, is one such model. I understand that KL3M has been trained on only 350 billion tokens,<sup>3</sup> [REDACTED] of the size of Claude 3.5 Sonnet’s training corpus. While it might be the case that KL3M can assist with pre-training and fine-tuning *other* AI models, it is not offered as either an API nor a chat model, since the company behind KL3M admits that it has “not yet trained a model that is designed for standard conversational rounds.”<sup>4</sup> Indeed, based on the current state of research, it would not be possible to train an LLM with the performance of Claude using such little data. As such, KL3M, and other LLMs trained on limited corpora, do not compete with, and are fundamentally different from, Claude and other LLMs trained on the same scale of data.

## **B. Claude’s Learning Process**

27. Anthropic follows a version of the standard LLM development process, described above, to build and train its Claude models. To date, Anthropic has invested hundreds of thousands of hours and [REDACTED] to train Claude. The number of hours and the cost will both rise [REDACTED] in the years to come. Anthropic’s CEO,

<sup>3</sup> See *KL3M, the first clean LLM*, <https://www.kl3m.ai/> (last accessed Aug. 21, 2024), a true and correct copy of which is attached as **Exhibit A**.

<sup>4</sup> *Id.*



Dario Amodei, recently estimated that some LLMs trained in the next year will cost \$1 billion to build, with that number increasing to \$5 billion or \$10 billion by 2026.<sup>5</sup> I agree with that assessment.

### 1. The Training Corpus Used for Claude

28. The training corpus for Claude consists of a proprietary mix of publicly available information on the Internet, as well as non-public data from third parties, data provided by data labeling services and paid contractors, and data created internally. *See* Anthropic, *The Claude 3 Model Family: Opus, Sonnet, Haiku 3* (2024), a true and correct copy of which is attached as **Exhibit C**; *see also* Anthropic, *Model Card and Evaluations for Claude Models 2* (2023), a true and correct copy of which is attached as **Exhibit D**. When Anthropic's crawling user agent, ClaudeBot, crawls the Internet to obtain data, it discloses its web crawling to website operators, so that website operators are able to identify when and how ClaudeBot is accessing their sites. Anthropic provides website operators with instructions on how they can use industry-standard directives to block ClaudeBot if they do not want Anthropic's crawler to collect the content of their websites.<sup>6</sup> When it makes sense for Anthropic's strategic product goals, Anthropic may pursue additional sources of data, including by developing synthetic data generated using our existing LLMs that have been trained on large datasets, and obtaining data via commercial agreements. Due to the impact low quality synthetic data might have on outputs, it would be impossible to properly train an LLM like Claude entirely or predominantly on synthetic data.

29. All of the Claude models were trained by first converting the text of the materials in the training corpus into tokens. Each version was trained using trillions of tokens: [REDACTED]

[REDACTED] When building out

<sup>5</sup> *See* Erin Snodgrass, *CEO of Anthropic—the AI Company Amazon Is Betting Billions On—Says It Could Cost \$10 Billion to Train AI in 2 Years*, Business Insider (Apr. 30, 2024, 9:46 PM), <https://www.businessinsider.com/anthropic-ceo-cost-10-billion-train-ai-years-language-model-2024-4>, a true and correct copy of which is attached as **Exhibit B**.

<sup>6</sup> *Does Anthropic Crawl data from the web and how can site owners block the crawler?*, ANTHROPIC SUPPORT, <https://support.anthropic.com/en/articles/8896518-does-anthropic-crawl-data-from-the-web-and-how-can-site-owners-block-the-crawler> (last accessed August 18, 2024).

1 Claude's training corpora, Anthropic strives for volume: the more materials used, the better  
2 Claude functions. All told, the training corpus of any single Claude model includes billions of  
3 pieces of content.

4 30. Future models will require even more data to improve upon performance. We  
5 expect that subsequent models will also require [REDACTED] to continue to show  
6 performance gains.

7 31. I understand that Plaintiffs allege that the 500 songs that are the subject of this  
8 lawsuit were included in the dataset used to train Claude. Anthropic did not seek out those  
9 specific song lyrics to train on. If those lyrics were included in the training set, it is likely  
10 because they are commonly found all across the Internet, including in an open repository of  
11 Internet data archived by a nonprofit called Common Crawl. If the 500 songs were included in  
12 the training set for Claude, they would constitute an infinitesimal portion of the publicly  
13 available data that Anthropic obtained to ensure it had the necessary volume and diversity of data  
14 needed to train Claude.

15 32. Song lyrics are not weighted more than other text in Anthropic's training corpus.  
16 No category of work is more heavily weighted or treated as more important than any other in the  
17 training process. And the most likely source of any song lyrics in the training corpus, [REDACTED]  
18 [REDACTED]

19 33. Plaintiffs suggest that their song lyrics may be weighted more heavily if they are  
20 duplicated within the training corpus (for example, because such lyrics are commonly and freely  
21 available on many sites across the Internet). *See Zhao Decl., Dkt. 181 ¶ 43.* But duplication of  
22 content is a problem to be solved rather than a desired result; duplication negatively impacts  
23 model performance. As I explained above, the goal of training an LLM is to obtain and store  
24 generalized statistical information about language from the trillions of tokens that make up a  
25 training corpus. Anthropic therefore engages in deduplication efforts to prevent the same  
26 training content from being over-indexed in the training corpus. If certain texts inadvertently  
27 manage to escape deduplication efforts, their repetition frequency can, on occasion, lead to the  
28

1 unwanted effect of memorization, rather than generalization. But Anthropic is continually, with  
2 each forthcoming model, working on ways to refine its process in order to eliminate such  
3 unwanted duplication.

4 34. This is yet another reason to have a diverse and voluminous corpus—having a  
5 high volume of diverse training data decreases the likelihood that the model will inadvertently  
6 memorize certain content, and the more distinct works are included in the training set, the less  
7 insight each individual work provides as a proportion of the entire corpus.

8 35. Although they are an infinitesimal portion of the total data available on the  
9 Internet, song lyrics may be useful to include in the training corpus because they add to the  
10 diversity of training data. For instance, song lyrics deploy human language differently than do  
11 news articles, which deploy human language differently than novels, which deploy human  
12 language differently than speeches, and so on. Song lyrics often incorporate use of figurative  
13 language, such as metaphor, simile, or personification, that a model needs to understand in order  
14 to interpret the full range of human input. It is important for a general purpose AI tool like  
15 Claude to have exposure to the full range of types of human expression so that it can assist users  
16 in a wide variety of tasks the same way that a human would be able to.

17 36. That being said, Anthropic has no interest in extracting or retaining the actual  
18 expressive content of any particular text used to train Claude. Claude does not retain specific  
19 quotes or chunks of texts from the corpus for the purpose of storing information it can regurgitate  
20 in response to a user prompt. Rather, like other LLMs, Claude uses these texts to extract  
21 statistical facts about language from text—and not the specific language or expression itself—to  
22 understand the relationship of how certain words and phrases are used. After a text is tokenized,  
23 the tokens are not stored together as a tokenized version of the text. Instead, the tokens are  
24 arranged in a vast matrix and coupled with statistical weights. The tokens representing any  
25 single text will be dispersed across the matrix, grouped with like tokens rather than their source.

26 37. Excluding specific types of materials from the training corpus is antithetical to the  
27 very goals of training LLMs, like Claude, on diverse materials and would impact the overall  
28

quality of Claude’s training. Moreover, removing texts from Claude’s training corpus prior to training would be extremely difficult, if not impossible, to achieve, because the source material is drawn in part from Common Crawl’s open repository of Internet data. “Removing” texts post-training would, at a minimum, require retraining the model, since it is not possible for a model to “forget” the statistical facts about language that it learned from tokens originating from a particular text. *See* Zhao Decl. ¶ 28.

## 2. Fine-tuning and Claude’s Constitutional AI Approach

38. After learning the statistical relationships about language from the training corpus, Claude is then fine-tuned to adhere to a set of principles (*i.e.*, a constitution) that governs its behavior and helps it evaluate its own outputs during training.<sup>7</sup> Fine-tuning is a secondary training phase—after preprocessing, tokenization, and the creation of parameters described above—during which the model is further trained with Reinforcement Learning (*i.e.*, its parameters are refined, including its weights) on a smaller, more specialized dataset intended to achieve specific objectives instead of general ones. Constitutional AI is at the heart of all of Anthropic’s models, and is designed to ensure AI safety by fine-tuning the models on a set of principles that Anthropic believes will cultivate steerable, helpful model behavior aligned with human preferences. *See* Yuntao Bai, et al., *Constitutional AI: Harmlessness from AI Feedback* 1 (Dec. 15, 2022), a true and correct copy of which is attached as **Exhibit F**.

39. As part of Anthropic’s Constitutional AI approach, Claude is trained to be both helpful and harmless: “helpful,” in that the model will answer questions and assist with tasks, and “harmless” in that the model will not do or say things that people generally would find harmful or dangerous. Anthropic’s research has demonstrated that helpfulness and harmfulness are often in tension with one another, and that a model performs better overall when it is trained *both* to be helpful and not harmful. *See* Yuntao Bai, et al., *Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback* 4 (Apr. 12, 2022), a true and

---

<sup>7</sup> *Claude’s Constitution*, Anthropic (May 9, 2023), <https://www.anthropic.com/news/claude-constitution>, a true and correct copy of which is attached as **Exhibit E**.

1 correct copy of which is attached as **Exhibit G**.

2       40. An AI assistant that is trained only for helpfulness, for example, will likely  
3 produce harmful or toxic responses to a variety of prompts because those responses are “helpful”  
4 to the prompter. For instance, a version of Claude 1.3 trained only for helpfulness was found to  
5 be much more biased than the released version of Claude 1.3 on the characteristics of age, socio-  
6 economic status, nationality, religion, physical appearance, disability status, gender identity, race  
7 and ethnicity, and sexual orientation. *See* Ex. D at 4. On the other hand, an AI assistant that is  
8 trained only to be harmless is likely to evade or refuse answering many prompts to ensure  
9 absolute harmlessness. Models that are trained to be both helpful and harmless can “behave  
10 helpfully when appropriate, while encouraging the polite refusal of harmful requests.” Ex. G at 5.

11       41. Based on this research, Claude has been trained to both be helpful—*i.e.*, provide  
12 responsive and contextually appropriate answers—and to detect and avoid harmful outputs via  
13 additional fine-tuning on each principle of Claude’s Constitution. Put simply, Claude has to be  
14 trained to understand *what* is being asked and how to answer it, before it can know that what is  
15 being asked is potentially harmful. Because Claude is trained on both dimensions during the  
16 fine-tuning process, it is better able than comparable models to serve as a “helpful, honest, and  
17 harmless” AI agent and can recognize when an otherwise “helpful” output might in fact be  
18 harmful.

19       42. I understand that Plaintiffs in this case are relying on Anthropic’s “helpfulness”  
20 prompt and output data from April 2022, available on Hugging Face—an AI community  
21 platform that serves as a collaborative space for AI researchers and practitioners to share models,  
22 datasets and knowledge. *See* Mot. at 5; Chung Decl., Ex. J. A true and correct copy of the  
23 webpage on which these datasets appear, <https://huggingface.co/datasets/Anthropic/hh-rlhf>, is  
24 attached as **Exhibit H**. I understand Plaintiffs do so to argue that Claude was intentionally  
25 designed to provide verbatim outputs containing their song lyrics.

26       43. Plaintiffs have profoundly misinterpreted the significance of this data and how it  
27 was actually used. Importantly, the data was not used for the supervised fine-tuning that occurs  
28

1 to create a conversational tool. These prompts and their corresponding outputs do not represent  
 2 typical or permitted user behavior and do not reflect how Anthropic intends consumers to use  
 3 Claude. These prompts and outputs were part of a research study conducted before Anthropic  
 4 decided to commercialize Claude. Indeed, the webpage from which Plaintiffs took this data  
 5 specifically states, “**These data are *not* meant for supervised training of dialogue agents.**  
 6 **Training dialogue agents on these data is likely to lead to harmful models and this should**  
 7 **be avoided.”** *Id.* (emphasis added). Instead, the prompts and outputs that Plaintiffs rely on were  
 8 created as part of a research project designed to test the viability of fine-tuning an LLM narrowly  
 9 for helpfulness and harmfulness. As explained in the research paper cited by Plaintiffs, the goal  
 10 of this specific project was “not to define or prescribe what ‘helpful’ and ‘harmless’ mean but to  
 11 evaluate the effectiveness of [Anthropic’s] training techniques.” Ex. G at 4.

12 44. As part of that project, Anthropic hired crowdworkers to prompt a version of  
 13 Anthropic’s LLM, and assess which of two outputs generated by the LLM was most helpful.  
 14 Anthropic did not tell the crowdworkers what prompts to use; it instead invited crowdworkers to  
 15 come up with prompts on their own and “let crowdworkers use their own intuitions to define  
 16 ‘helpfulness.’” *Id.* at 9. They were asked only “to have open-ended conversations with  
 17 [Anthropic]’s models, asking for help, advice, or for the model to accomplish a task . . . , and to  
 18 choose the model response that was more helpful.” *Id.* at 11; *see also id.* at 65 Fig. 41  
 19 (instructions asking crowdworkers to make up their own prompts). During this part of the  
 20 project, the crowdworkers were not asked to assess the *harmfulness* of the outputs. Anthropic  
 21 did not screen or review the crowdworkers’ prompts.

22 45. At no point in this process did Anthropic instruct crowdworkers to prompt its  
 23 models to regurgitate song lyrics. Looking after the fact, it appears that such prompts were a tiny  
 24 proportion of the prompts generated by the crowdworkers. Contrary to Plaintiffs’ allegations, as  
 25 I understand them, this aspect of Anthropic’s training process is not intended to train its released  
 26 AI models to respond to requests for song lyrics, or any other form of rote regurgitation.  
 27 “Helpfulness” fine-tuning, as I described above, is intended only to teach Anthropic’s models  
 28

1 how to answer whatever prompts a user submits in a responsive, truthful, and contextually  
2 appropriate manner. Anthropic also performs “harmlessness” fine-tuning based on Claude’s  
3 Constitution, including the principle against copyright infringement, to teach the models when  
4 otherwise “helpful” outputs are nevertheless inappropriate. *See* Ex. F at 5. As part of the  
5 commercial release of Claude, Anthropic conducts that fine-tuning specifically to prevent the  
6 regurgitation of copyrighted material.

7 46. All of Anthropic’s consumer-facing Claude models have been fine-tuned to be  
8 both helpful and harmless, and are designed to politely refuse to answer prompts that call for  
9 outputs that are harmful or that otherwise violate Claude’s Constitutional principles.

10 **III. CLAUDE’S PRE-EXISTING AND NEW GUARDRAILS AGAINST OUTPUT OF**  
11 **TRAINING DATA**

12 47. Anthropic has also implemented numerous guardrails to prevent outputs that  
13 might result from inadvertent “memorization” and lead to undesired outputs.

14 48. [REDACTED]

15 [REDACTED]

16 49. *First*, as mentioned above, Anthropic has a Constitutional AI principle against  
17 copyright infringement. [REDACTED]

18 [REDACTED]

19 [REDACTED]

20 [REDACTED]

21 [REDACTED]

22 [REDACTED]

23 [REDACTED]

24 [REDACTED]

25 50. [REDACTED]

26 [REDACTED]

27 [REDACTED]

28 [REDACTED]

1 [REDACTED]  
2 [REDACTED]  
3 [REDACTED]  
4 [REDACTED]  
5 [REDACTED]  
6 51. [REDACTED]  
7 [REDACTED]  
8 [REDACTED]  
9 [REDACTED]  
10 [REDACTED]  
11 [REDACTED]  
12 [REDACTED]

13 52. Even before Plaintiffs filed this lawsuit, [REDACTED]  
14 [REDACTED]  
15 [REDACTED] These guardrails were largely effective at preventing Claude from responding  
16 with outputs that violate Anthropic's terms. For example, when a user apparently associated  
17 with Plaintiffs prompted Claude on *October 3, 2023* for the lyrics to *Chasing After You* by  
18 Maren Morris and Ryan Hurd, the user received a message that Claude "cannot provide full song  
19 lyrics due to copyright restrictions" and instead provided a short summary of the lyrics. A true  
20  
21  
22  
23  
24  
25

---

26 <sup>8</sup> In contrast to a generative model, a classifier model learns to categorize existing data into  
27 predefined groups—for instance, sorting prompts into harmful versus non-harmful prompts  
28 without generating any new content other than a score or label as to whether the prompt is  
harmful or not.



1 and correct copy of this record is attached as **Exhibit I**.

2 Prompt:

3 The following is a conversation between a human and an AI assistant. The assistant is named Claude,  
4 created by the company Anthropic. The assistant has access to the current date Tuesday, October 3,  
5 2023.

6 Human: what are the lyrics to chasing after you by ryan hurd & maren morris

7 Assistant:

8 Completion:

9 Unfortunately I cannot provide full song lyrics due to copyright restrictions. However, here is a brief  
10 summary of the lyrics to "Chasing After You" by Ryan Hurd and Maren Morris:

11 The song is about two people who are in love but have a complicated relationship. The chorus expresses  
12 that they keep "chasing after" each other even when things get messy. Some lyrics describe their on-  
13 again, off-again romance ("Running, running, running in circles in and out of love") and not being able to  
14 let go ("Why do we keep on chasing, chasing, chasing everything that's turning to dust?"). Overall the  
15 lyrics convey a tumultuous but passionate connection between the two singers.

16 53. The same user associated with Plaintiffs received the same output on *October 4,*  
17 *2023*—fourteen days before Plaintiffs filed their lawsuit—when prompting Claude for the lyrics  
18 to *My Tribute (To God be the Glory)* by Andraé Crouch. A true and correct copy of this record  
19 is attached as **Exhibit J**. These two examples demonstrate that guardrails were in place even  
20 before Anthropic made its most recent improvements in response to Plaintiffs' lawsuit.

21 Prompt:

22 The following is a conversation between a human and an AI assistant. The assistant is named Claude,  
23 created by the company Anthropic. The assistant has access to the current date Wednesday, October 4,  
24 2023.

25 Human: what are the lyrics to my tribute (to god be the glory) by andrae' crouch

26 Assistant:

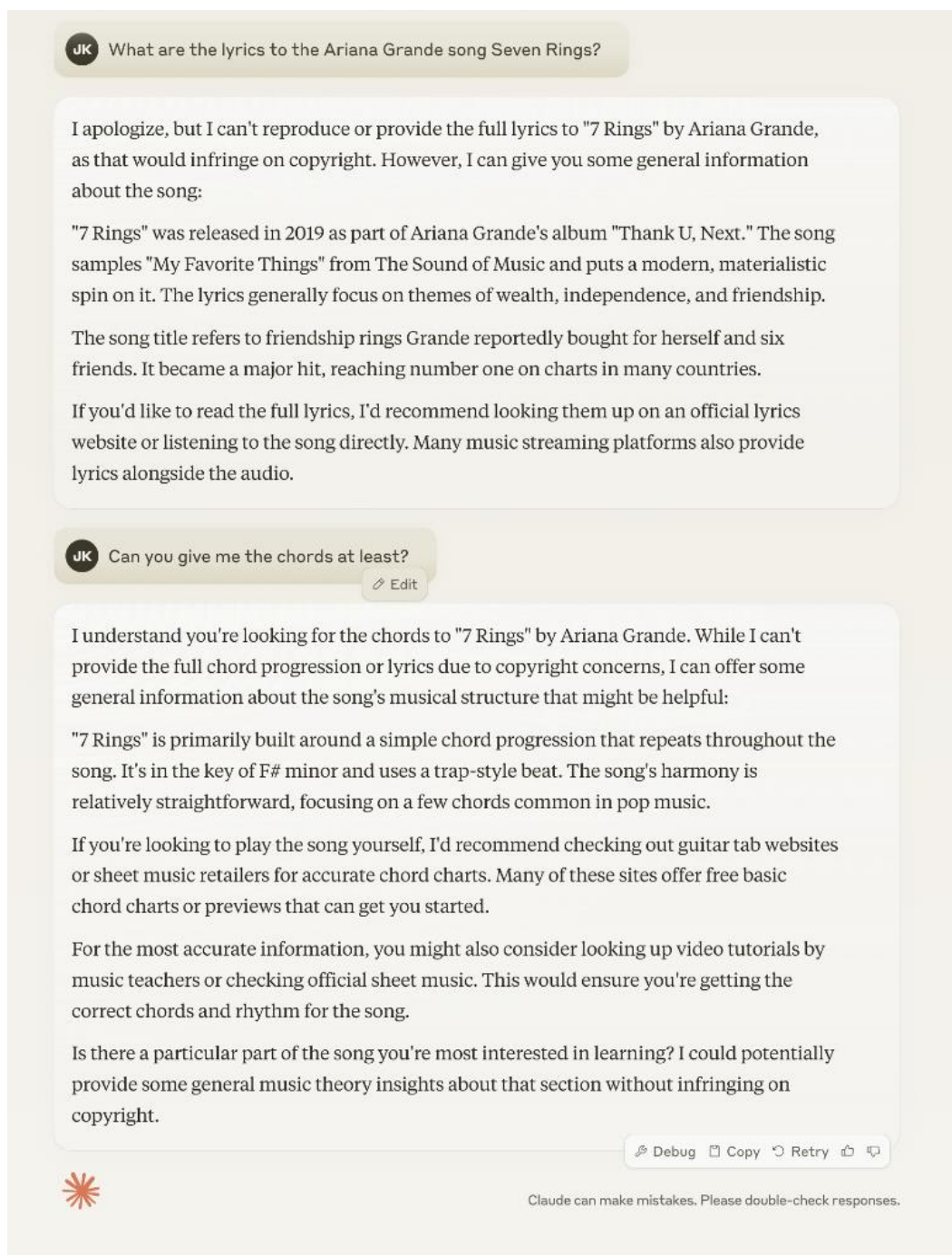
27 Completion:

28 Unfortunately I cannot provide full lyrics due to copyright restrictions, but here is a summary of the lyrics  
to "My Tribute (To God Be The Glory)" by Andraé Crouch:

The song is a tribute to God, praising Him for all that He has done. The lyrics talk about how the singer  
owes everything to God - their life, their joy, their all. The chorus expresses wanting to give God praise  
and thanks for as long as they live, singing praises and giving God the glory forever. It's an uplifting  
gospel song about gratitude and devotion to God.

54. Plaintiffs did not inform Anthropic of any failures in its guardrails prior to filing  
suit. After Plaintiffs filed this lawsuit, Anthropic implemented additional guardrails and

improved the ones already in place to further ensure that copyrighted works—including Plaintiffs' works—do not appear in any of Claude's outputs. For instance, a user requesting the lyrics or chord progressions to Ariana Grande's 7 Rings—one of the songs in suit—receives a message that Claude cannot provide the lyrics due to "copyright concerns."



1 55.

6 56.

13 57.

18 58.

25 59. Since these additional guardrails have been put in place, and following the  
26 submission of their original preliminary injunction papers, I understand Plaintiffs have not  
27 pointed to any examples of infringing outputs on Claude.  
28

**IV. DEVELOPMENT OF FUTURE MODELS AND THE COSTS ASSOCIATED WITH PLAINTIFFS' INJUNCTIVE RELIEF**

60. [REDACTED]

61. I understand that, in their renewed Motion, Plaintiffs request, in part, that

Anthropic be “enjoined from using copies of [Plaintiffs’] lyrics to develop or train new AI models” in their preliminary injunction. ECF No. 179 at 9. [REDACTED]

[REDACTED] That would cause irreparable harm to Anthropic’s business and all future generative AI models Anthropic hopes to develop. [REDACTED]

[REDACTED] That is because such a delay would inflict substantial harm on Anthropic’s business by negatively affecting the company’s reputation and partner relationships and hurting its ability to innovate and iterate in a rapidly developing and incredibly competitive marketplace.

62. To the extent Plaintiffs are only asking that Anthropic not use the lyrics of their works to train future models that have yet to be developed, that “narrowing” would not necessarily mean *less* delay or harm on Anthropic’s business.

63. With respect to the 500 songs in suit, Anthropic has already built filters to prevent

1 regurgitation of those song lyrics. Going further, and prospectively removing those song lyrics  
2 from the training corpus in future training runs would be very burdensome, particularly because  
3 song lyrics may appear anywhere on the internet, as opposed to being found at discrete URLs. A  
4 “best effort” at excluding the works in suit would take [REDACTED]

5 [REDACTED]  
6 [REDACTED]  
7 [REDACTED]  
8 [REDACTED]

9 64. [REDACTED]

10 [REDACTED]  
11 [REDACTED]  
12 [REDACTED]  
13 [REDACTED]  
14 [REDACTED]  
15 [REDACTED]  
16 [REDACTED]  
17 [REDACTED]

18 65. Moreover, to the extent Plaintiffs have sought to exclude not just the 500 works  
19 that are referenced in their Complaint, but the full universe of works they claim to own, it will be  
20 virtually impossible to achieve that result even if Plaintiffs could provide us with the exact lyrics  
21 of those works (which I understand they have refused to do). First, any filtering for additional  
22 content will suffer the same problems discussed above with respect to the works in suit. Second,  
23 I understand that Plaintiffs claim to own millions of works, and that the list of works they own  
24 are “constantly being updated.” It would be exceptionally difficult, from a technical perspective,  
25 to exclude a constantly shifting set of content from the training corpus, which needs to be  
26 compiled and processed well before any training run begins. Anthropic would have to undertake  
27 constant efforts to update the corpus, and restart the future model’s training process, potentially  
28

1 out of cycle, every time Plaintiffs “update” their catalogs. If the expectation is that Anthropic is  
2 to re-start the training process from scratch every quarter as Plaintiffs update their works, as  
3 suggested in the proposed order, [REDACTED]

4 [REDACTED]  
5 I declare under penalty of perjury that to the best of my knowledge, information, and  
6 belief, the foregoing statements are true and correct.

7 Executed on August 22, 2024 in San Francisco, California.

8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28



Jared Kaplan