

1 Rom Bar-Nissim (SBN 293356)
Rom@HeahBarNissim.com
2 HEAH BAR-NISSIM LLP
1801 Century Park East, Suite 2400
3 Los Angeles, CA 90067
Telephone: (310) 432-2836

4 Jarrett Lee Ellzey (*pro hac vice application*
5 *forthcoming*)
Tom Kherkher (*pro hac vice application*
6 *forthcoming*)
Leigh S. Montgomery (*pro hac vice application*
7 *forthcoming*)
ELLZEY KHERKHER SANFORD
8 MONTGOMERY LLP
4200 Montrose Street, Suite 200
9 Houston, TX 77006
Jellzey@EKSM.com
10 TKherkher@EKSM.com
LMontgomery@EKSM.com

11
12 *Attorneys for Plaintiffs and the Proposed Class*

13
14 **UNITED STATES DISTRICT COURT**
15 **NORTHERN DISTRICT OF CALIFORNIA**
16 **SAN FRANCISCO DIVISION**
17

18 TED ENTERTAINMENT, INC., MATT
19 FISHER, and GOLFHOLICS, INC., each
20 individually and on behalf of all others
similarly situated,

21 Plaintiffs,

22 v.

23 OPENAI, INC., OPENAI GP LLC, OPENAI
24 OPCO LLC, OPENAI GLOBAL LLC, OAI
CORPORATION, LLC, OpenAI Group
25 PBC, and OPENAI HOLDINGS, LLC

26 Defendants
27
28

Case No.: 3:26-cv-2935

CLASS ACTION COMPLAINT

DEMAND FOR JURY TRIAL

1 Plaintiffs Ted Entertainment, Inc. (“TEI”), Matt Fisher, and Golfholics, Inc. (collectively,
2 where appropriate “Plaintiffs”), each individually and on behalf of all others similarly situated (the
3 “Class,” as defined below), by and through their undersigned counsel, file this Complaint against
4 OpenAI, Inc.; OpenAI GP LLC; OpenAI OpCo LLC; OAI Corporation, LLC; OpenAI Holdings,
5 LLC; and OpenAI Group PBC (collectively, “OpenAI” or “Defendant”), and alleges as follows:

6 **NATURE OF THE ACTION**

7 1. This lawsuit arises from Defendant unlawfully circumventing technological
8 protection measures to access and scrape millions of copyrighted videos from the online video
9 viewing platform, YouTube, in order to feed, train, improve, and commercialize Defendant’s large-
10 scale generative artificial intelligence (“AI”) model named “Sora.”

11 2. Sora is a text-to-video AI model conceived and marketed as OpenAI's dedicated
12 video generation platform. Sora was significant enough that Defendant launched it as its own
13 standalone application, the Sora app, separate from and in addition to its integration within
14 ChatGPT.

15 3. YouTube allows the public to view audiovisual works only through controlled
16 streaming, keeping its underlying video files locked away behind careful access restrictions.
17 Defendant picked this lock by deploying scraping tools designed to evade and circumvent
18 YouTube's access restrictions. Defendant’s circumvention of YouTube's technological protection
19 measures (“TPMs”) was systematic. On information and belief, Defendant used automated video-
20 downloading programs combined with virtual machines that rotated IP addresses to avoid detection
21 and blocking, enabling the mass unauthorized access and extraction of videos at the scale necessary
22 to train Sora. Defendant then used Plaintiffs' and Class Members' intellectual property for their own
23 commercial gain in developing and commercializing Sora. In doing so, Defendant has violated the
24 law and YouTube's Terms of Service, which were intended to protect Plaintiffs and others similarly
25 situated.

26 4. By accessing, scraping, and downloading those files, Defendant deliberately
27 circumvented YouTube's access controls to obtain the training data necessary to build and
28 commercialize Sora.

1 5. Plaintiffs and the Class Members whom Plaintiffs seek to represent are content
2 creators who upload their audiovisual content to YouTube. In doing so, the content creators are
3 authorizing and instructing YouTube to provide protection to the video content through YouTube's
4 anti-circumvention software and Terms of Service. In fact, YouTube's anti-circumvention software
5 and protective Terms of Service are a driving factor behind many content creators' decision to
6 upload their video content to YouTube.

7 6. Rather than seek permission or pay a fair price for the audiovisual content hosted on
8 YouTube, Defendant harvested content creators' protected and copyrighted videos for commercial
9 use and at scale without consent or compensation to the content creators.

10 7. Defendant's actions were not only unlawful, but an unconscionable attack on the
11 community of content creators whose content is used to fuel the multi-trillion-dollar generative AI
12 industry without any compensation.

13 8. Content creators such as Plaintiffs and the Class Members will never be able to claw
14 back the intellectual property unlawfully copied and used by Defendant to train its generative AI
15 models. Once AI ingests content, that content is stored in its neural network and not capable of
16 deletion or retraction. Defendant's actions constitute abuse and exploitation of content creators'
17 work for Defendant's profit.

18 9. Most YouTube videos are not registered with the U.S. Copyright Office. That lack
19 of registration, however, does not render them valueless or leave them unprotected. Content creators
20 invest time, skill, and resources into producing their works, and they rely on YouTube's TPMs to
21 safeguard their files from unauthorized access. Because copyright registration is not a prerequisite
22 for protection against unlawful circumvention of access controls, this claim is especially critical
23 where Defendant's misconduct lies in breaking through access barriers that prevent anyone from
24 obtaining the underlying files in the first place.

25 10. It is also critical to protect a flourishing internet ecosystem. In a world where
26 Defendant and others can circumvent technological protections to exploit copyrighted works
27 without authorization with impunity, creators will be less likely to make their creations available on
28 YouTube and other similar platforms, for fear of losing all control of them. The world will be poorer

1 for it.

2 11. The Class Members, including Plaintiffs, are content creators who uploaded their
3 audiovisual content to YouTube’s video platform. On information and belief, the audiovisual
4 content of Class Members was among the video content utilized by Defendant to train Sora and
5 related generative AI models.

6 12. An essential component of Defendant’s business model—powering AI features and
7 services with large-scale training data—includes the mass acquisition and ingestion of creators’
8 videos scraped from YouTube.

9 13. Defendant has profited substantially from its infringement of Plaintiffs’ and Class
10 Members’ video content through Defendant’s training of its generative AI products. Defendant’s
11 massive financial success would not have been possible without the video content created by
12 Plaintiffs and Class Members, which was intended for streaming on YouTube.

13 14. Plaintiffs bring this class action on behalf of themselves and on behalf of a
14 nationwide class of YouTube creators whose works were scraped, ingested, and trained on without
15 authorization, seeking statutory damages, injunctive relief, restitution, and all other remedies
16 allowed by law pursuant to the Digital Millennium Copyright Act, 17 U.S.C.A. § 1201(a) seeking
17 an injunction and damages commensurate with the scope of Defendant’s massive and ongoing
18 infringement. More particularly, Defendant’s conduct violates the provisions of the DMCA
19 regarding anti-circumvention (§1201) by bypassing technological protection measures that control
20 access to YouTube videos.

21 **THE PARTIES**

22 15. Plaintiff Ted Entertainment, Inc. (“TEI”), is an independent California based media
23 company and content creator with over 5,800 original videos on YouTube, with a combined total of
24 over 4 billion views. TEI has amassed a substantial following on YouTube with over 2.6 million
25 subscribers to its channels. TEI owns and operates the channels “h3h3 Productions” and “H3 Podcast
26 Highlights.” H3H3 Productions appears with 146 videos in HD-VILA-100M, 83 videos in HD-VG-
27 130M, 155 videos in Panda-70M, and 24 videos in HowTo100M. H3 Podcast Highlights appears
28 with 285 videos in HD-VILA-100M, 209 videos in HD-VG-130M, 283 videos in Panda-70M, and

1 1 video in HowTo100M. Plaintiff invested significant resources – time and money – into producing
2 and publishing this video content and bringing awareness to its content.

3 16. Plaintiff TEI is the creator of the videos identified in Exhibit A, all of which were
4 uploaded exclusively to YouTube by Plaintiff. Plaintiff derives value from the works identified in
5 Exhibit A through viewership, advertising, sponsorships, licensing, and related monetization.

6 17. Plaintiff TEI and its owners – Ethan and Hila Klein – are longtime champions of the
7 rights of YouTube creators. The Kleins helped define what constitutes fair use reaction videos and
8 the good faith belief standard for DMCA counternotifications. TEI has championed online free
9 speech and is currently helping define what reaction content does not constitute fair use to ensure
10 YouTube content creators can enjoy the fruits of their labor.

11 18. Plaintiff Matt Fisher is an individual and resident of the State of California. He is a
12 golf content creator who posts his original videos on YouTube, many of which are instructional.
13 His channel is “MrShortGame Golf” on the YouTube platform. He has over 500,000 subscribers
14 and hundreds of millions of views. Plaintiff has invested substantial time and money into bringing
15 awareness around his content. MrShort Game appears with 2 videos in HD-VILA-100M, 1 video in
16 HD-VG-130M, 8 videos in Panda-70M, and 4 videos in HowTo100M.

17 19. Plaintiff Matt Fisher is the creator of the videos identified in Exhibit B, all of which
18 were uploaded exclusively to YouTube by Plaintiff and all of which were included in datasets used
19 by Defendant to train its generative AI models. Plaintiff derives value from the works identified in
20 Exhibit B through viewership, advertising, sponsorships, licensing, and related monetization.

21 20. Plaintiff Golfholics is a corporate entity organized pursuant to the laws of the State
22 of California. It is a golf content channel that posted its original videos on YouTube. The channel is
23 “Golfholics” on the YouTube platform. It has over 130,000 subscribers and millions of views.
24 Plaintiff invested substantial time and money into bringing awareness around its content. Golfholics
25 appears with 62 videos in HD-VILA-100M, 62 videos in Panda-70M, and 37 video in HD-VG-
26 130M.

27 21. Plaintiff Golfholics is the creator of the videos identified in Exhibit C, all of which
28 were uploaded exclusively to YouTube by Plaintiff and all of which were included in datasets used

1 by Defendant to train its generative AI models. Plaintiff derives value from the works identified in
2 Exhibit C through viewership, advertising, sponsorships, licensing, and related monetization.

3 22. Defendants are closely related Delaware entities that together operate as OpenAI.
4 Defendant OpenAI Inc. is a Delaware nonprofit corporation with a principal place of business
5 located in San Francisco, California. OpenAI Inc. indirectly owns and controls all other OpenAI
6 entities, including all other Defendants, and has been directly involved in carrying out the large-
7 scale 17 U.S.C. §1201(a) anti-circumvention violations alleged in this Complaint.

8 23. Defendant OpenAI GP, LLC is a Delaware limited liability company with a
9 principal place of business located in San Francisco, California. OpenAI GP LLC is wholly owned
10 and controlled by OpenAI Inc. OpenAI GP, LLC has been directly involved in carrying out the
11 large-scale 17 U.S.C. §1201(a) anti-circumvention violations alleged in this Complaint through its
12 direction and control of OpenAI LP and OpenAI Global LLC.

13 24. Defendant OpenAI Global, LLC is a Delaware limited liability company with a
14 principal place of business located in San Francisco, California. OpenAI Global LLC owns, sells,
15 licenses, and monetizes a number of OpenAI's offerings, including ChatGPT, ChatGPT Enterprise,
16 and OpenAI's API tools, all of which were built on OpenAI's large-scale 17 U.S.C. §1201(a) anti-
17 circumvention violations alleged in this Complaint. Upon information and belief, OpenAI Global,
18 LLC is owned and controlled by OpenAI Inc.

19 25. Defendant OpenAI OpCo LLC is a Delaware limited liability company with a
20 principal place of business located in San Francisco, California. OpenAI OpCo LLC is a wholly
21 owned subsidiary of OpenAI Inc. and has facilitated and directed OpenAI's largescale 17 U.S.C.
22 §1201(a) anti-circumvention violations alleged in this Complaint through its management and
23 direction of OpenAI Global, LLC.

24 26. Defendant OAI Corporation, LLC is a Delaware limited liability company with a
25 principal place of business located in San Francisco, California. OAI Corporation, LLC's sole
26 member is OpenAI Holdings, LLC/OpenAI Group PBC. OAI Corporation, LLC was and is involved
27 in the large-scale 17 U.S.C. §1201(a) anti-circumvention violations alleged in this Complaint
28 through its ownership, control, and direction of OpenAI Global LLC and OpenAI LLC.

1 **FACTUAL BACKGROUND**

2 **A. YouTube’s Terms of Service Prohibit Scraping, Downloading, and**
3 **Unauthorized Access to Video Files, and Employs Technological Measures**
4 **Consistent With This Prohibition**

5 33. As noted above, Plaintiffs and the Class Members are content creators of audiovisual
6 content.

7 34. Plaintiffs and the Class Members create original video content that has been
8 uploaded onto YouTube’s video sharing platform.

9 35. YouTube’s “users”—the individuals who view the digital content available on
10 YouTube—can watch and listen to videos for free on YouTube’s advertising-supported service, but
11 YouTube does not give users access to or allow downloading of the digital files underlying the
12 content viewed by the user through YouTube’s authorized playback mechanism.

13 36. YouTube’s Terms of Service expressly prohibit scraping, unauthorized
14 downloading, bulk extraction, or other forms of data mining of audiovisual content except through
15 expressly permitted features or licensed APIs. These contractual restrictions operate together with
16 YouTube’s TPMs to prevent unlicensed access to creators’ videos.

17 37. According to YouTube’s Terms of Service, content creators such as Plaintiffs who
18 upload content onto YouTube grant license to YouTube for certain uses as well as to other users of
19 YouTube to access content through YouTube’s services; however, the license makes clear that it
20 “does not grant any rights or permissions for a user to make use of [the] Content independent of the
21 Service.”¹

22 38. This language confirms that end users are not given access to the file itself, only the
23 ability to view (*i.e.*, stream) through YouTube’s controlled environment. YouTube’s Terms of
24 Service reflect YouTube’s intent to restrict access to the digital files underlying the videos that
25 YouTube’s users are allowed to stream through YouTube’s platform.

26 39. Streaming through YouTube and downloading permanent copies provide the user
27 with different value propositions—watching and listening for free but seeing ads, versus possessing

28 ¹ “Terms of Service,” YouTube, <https://www.youtube.com/t/terms> (last viewed March 11, 2026).

1 a permanent digital copy.

2 40. Accordingly, scraping or bulk downloading is not merely copying material already
3 provided; it is an act of unauthorized access to data files that YouTube affirmatively withholds from
4 public download.

5 41. In fact, during a Bloomberg interview about OpenAI scraping YouTube’s CEO,
6 Neal Mohan, confirmed that unauthorized scraping constitutes a violation of creators’ rights and
7 YouTube’s Terms of Service stating, **“From a creator’s perspective, when a creator uploads their
8 hard work to our platform, they have certain expectations. One of those expectations is that
9 the terms of service is going to be abided by,”** Mohan said, **“It does not allow for things like
10 transcripts or video bits to be downloaded, and that is a clear violation of our terms of service.
11 Those are the rules of the road in terms of content on our platform.”**²

12 42. The scraping and acquisition processes used by Defendant to circumvent YouTube’s
13 TPMs were inconsistent with, and in violation of, YouTube’s Terms of Service, which forbid
14 scraping and mass downloading of videos.³

15 43. To enforce its prohibitions on downloading content, YouTube does not provide a
16 downloading option that is readily available to users.

17 44. Although YouTube offers downloading options to subscribers who pay for
18 YouTube’s “Premium” plan, YouTube still employs TPMs to restrict access. Specifically, YouTube
19 prohibits all downloading audiovisual content except to the YouTube app. Further, the “download”
20 option only makes audiovisual content available for offline streaming—it does not provide the
21 Premium subscriber with access to the audiovisual files. Accordingly, the audiovisual files cannot
22 be transferred to any other device, but remain only for streaming on the app. Finally, the files are
23 available only for offline streaming for a limited amount of time, at which point they will become
24 available only for online streaming once again.

25 ² Davey Alba & Emily Chang, *YouTube Says OpenAI Training Sora With Its Videos Would Break*
26 *Rules*, Bloomberg (Apr. 4, 2024), <https://www.bloomberg.com/news/articles/2024-04-04/youtube-says-openai-training-sora-with-its-videos-would-break-the-rules> (video available at
27 https://www.youtube.com/watch?v=FBZ_BeChRg).

28 ³ *Supra*, note 1 (prohibiting “access, reproduce, download, distribute, transmit, broadcast, display, sell, license, alter, modify or otherwise use any part of the Service or any Content except: (a) as expressly authorized by the Service; or (b) with prior written permission from YouTube and, if applicable, the respective rights holders;”).

1 45. For users who do not have a “Premium” plan, the “download” option on YouTube’s
2 player page is non-functional.

3 46. To further enforce its prohibitions on downloading content, YouTube uses
4 technological processes and tools to detect and block access to files for unauthorized downloading.
5 For example, YouTube monitors downloading activity and may block IP addresses that make too
6 many download attempts in a specified period.

7 47. YouTube has employed a variety of TPMs to restrict access to Class Members’
8 audiovisual files. For example, YouTube does not provide users with a readily available
9 downloading option and it has implemented a number of tools, such as the Rolling Cipher and IP
10 Blocking, that operate to restrict users’ ability to access audiovisual material.

11 48. Content creators, including Plaintiffs and Class Members, rely on the TPMs and
12 YouTube’s Terms of Service in deciding to upload their audiovisual content to YouTube. Content
13 creators, including Plaintiffs and Class Members, expect that their works will not be copied at scale
14 without consent.

15 **B. YouTube’s Technological Protection Measures Function as Access Controls**
16 **Over Video Files.**

17 49. YouTube deploys multiple technological protection measures (“TPMs”) designed
18 to control, restrict, and monitor access to the underlying video files and to deter direct downloading
19 or bulk extraction outside YouTube’s controlled playback environment. These TPMs include,
20 among others: (1) an obfuscated signature system commonly referred to as a “rolling cipher,” (2)
21 IP-based blocking and rate limiting that restrict high-volume automated access, (3) short-lived,
22 session-bound streaming URLs, (4) CAPTCHA human-verification challenges triggered by
23 automated activity, and (5) proof-of-origin tokens that verify requests originating from authorized
24 client environments. Each TPM independently functions as a gatekeeping mechanism that must be
25 satisfied before the audiovisual file can be retrieved, and circumvention of any one of them enables
26 access to the file in a manner not authorized by YouTube or the content owners.

27 50. **TPM (1) - Rolling cipher (access control):** YouTube controls access to the
28 underlying media file by withholding a usable file location unless the requesting client can transform

1 an obfuscated signature parameter using proprietary logic embedded in the official YouTube player.
2 The playback data delivered to users contains a scrambled signature that must be processed to
3 produce a valid media request URL. Without performing this authorized transformation, the server
4 will not deliver the audiovisual file. Thus, the rolling cipher controls access by requiring application
5 of a specific computational process before the file can be obtained at all.

6 51. YouTube's rolling cipher encryption measure acts as a digital lock, controlling
7 access to content by protecting against unauthorized downloading of underlying media files.
8 YouTube maintains two different URLs for any given video: the page URL, visible to the user, is
9 for the webpage where the video playback occurs, and the file URL, not visible to the user, is for
10 the video file itself that is played within the page. The file URL is encrypted using a complex and
11 periodically changing algorithm – the rolling cipher – that is designed to impede external access to
12 the underlying YouTube files. YouTube's player software uses a decryption routine to authenticate
13 requests and deliver content only through approved interfaces. This TPM inhibits access to the
14 underlying audiovisual files for the purposes of any downloading, copying or distribution of the
15 audiovisual content. In other words, the rolling cipher controls access to copies of audiovisual
16 content uploaded to the platform and impedes an ordinary user from creating a permanent,
17 unrestricted download of audiovisual content made available on YouTube only for streaming and
18 restricted downloading.

19 52. The rolling cipher is a TPM measure within the meaning of 17 U.S.C. §1201(a)
20 because it “effectively controls access” to copyrighted works by preventing users access to the files
21 of video streams without first executing YouTube's proprietary decryption code.

22 53. **TPM (2) - IP blocking and rate limiting (access control):** YouTube monitors
23 network behavior and restricts access when excessive or abusive request patterns are detected.
24 Requests from an IP address that exceeds defined thresholds may be throttled, denied, or blocked
25 entirely, preventing further retrieval of audiovisual data from that source. This mechanism controls
26 access by limiting how frequently and at what scale a particular requester may obtain the file and by
27 cutting off access altogether when abuse is detected.

28 54. YouTube's infrastructure monitors the number and frequency of video requests

1 originating from individual IP addresses. This system serves an important function. Even if a user
2 were able to generate a valid media request for an individual video, YouTube’s infrastructure is
3 designed to prevent automated systems from retrieving large numbers of videos at scale.

4 55. Once YouTube detects abnormal volume of access to files, the system can block the
5 offending IP address from accessing YouTube’s servers. This prevents automated scraping systems
6 from accessing large volumes of copyrighted audiovisual content.

7 56. This IP-based monitoring and blocking system therefore also functions as a TPM
8 measure within the meaning of 17 U.S.C. §1201(a) because it “effectively controls access” to
9 copyrighted works by restricting automated access to YouTube’s underlying video files.

10 57. **TPM (3) - Session-bound, short-lived URLs (access control):** Even after a valid
11 media URL is generated, YouTube restricts access by issuing URLs that are temporary,
12 cryptographically signed, and tied to a particular playback session and client context. These URLs
13 include authorization parameters such as expiration timestamps and client identifiers. Once the
14 authorization window expires or the session ends, the server will refuse to deliver the audiovisual
15 file in response to that URL. Accordingly, possession of a previously valid link does not provide
16 continuing access to the file, and new authorization must be obtained to retrieve it. This mechanism
17 therefore controls access by limiting when and from where the file may be requested.

18 58. Because session-bound URLs expire automatically, any automated system
19 attempting to access videos outside ordinary playback cannot rely on a static link and must instead
20 repeatedly obtain fresh authorization parameters and regenerate valid media URLs. By
21 programmatically renewing expired credentials and initiating new authorized sessions at machine
22 speed, automated scraping tools can maintain continuous access to audiovisual files that would
23 otherwise become unavailable once the original authorization lapses. This conduct circumvents the
24 temporal and session-based restrictions imposed by YouTube and allows persistent retrieval of files
25 in a manner not available to ordinary users. Circumventing this measure constitutes circumvention
26 of a technological measure that effectively controls access to a copyrighted work in violation of 17
27 U.S.C. § 1201(a).

28 59. **TPM (4) - CAPTCHA challenges (access control):** When traffic patterns indicate

1 automated or suspicious activity, YouTube may require completion of a CAPTCHA challenge
2 before allowing further requests to proceed. Until the challenge is successfully completed, the
3 requesting client cannot obtain the data necessary to retrieve the audiovisual file. CAPTCHAs
4 therefore function as an access control by conditioning continued access to content on proof that the
5 requester is a human user rather than an automated system.

6 60. Large-scale scraping operations necessarily generate request patterns that would
7 ordinarily trigger such human-verification challenges. To continue accessing without interruption,
8 automated systems must be configured either to evade detection, distribute requests across multiple
9 machines, or otherwise prevent CAPTCHA enforcement from halting access. By avoiding or
10 neutralizing these verification gates, such systems obtain audiovisual data without demonstrating
11 the human interaction that YouTube requires for continued access. Defeating or bypassing this
12 human-verification mechanism constitutes circumvention of a technological measure that
13 effectively controls access to copyrighted works in violation of 17 U.S.C. § 1201(a).

14 61. **TPM (5) - Proof-of-origin tokens (access control):** YouTube further restricts
15 access by requiring that requests for video segments include cryptographic tokens demonstrating
16 that the request originates from an authorized client environment operating within the intended
17 playback context. These tokens are generated during playback and validated by YouTube's servers
18 before data is delivered. Requests lacking valid tokens, or presenting forged or replayed tokens, are
19 denied or degraded. This mechanism controls access by preventing unauthorized software from
20 directly requesting the audiovisual file unless it can present credentials proving it is an approved
21 client.

22 62. Automated tools operate outside YouTube's authorized player environment and
23 therefore must extract, replicate, or reuse the necessary request parameters in order to obtain access
24 to video data. By presenting such credentials outside their intended context, these tools can retrieve
25 audiovisual files without operating as approved clients. This conduct bypasses the origin-
26 verification function of the token system and enables direct access to media files that would
27 otherwise be delivered only to authorized playback software. Such bypassing constitutes
28 circumvention of a technological measure that effectively controls access to copyrighted works in

1 violation of 17 U.S.C. § 1201(a).

2 **C. Defendant Improperly Circumvented YouTube’s Technological Protection**
3 **Measures to Obtain Millions of YouTube Videos to Train Its Foundational AI**
4 **Model.**

5 63. As noted above, Defendant require significant amounts of data to feed, train,
6 improve, and commercialize Sora.

7 64. The full extent of Defendant’s unauthorized access and extraction of video content
8 to train Sora is not yet known and will be the subject of discovery. Defendant has compounded this
9 uncertainty by deliberately concealing the origins of Sora's training data behind the vague label of
10 “publicly available” data. When OpenAI CTO Mira Murati was directly asked what data was used
11 to train Sora, her answer was: “We used publicly available data...”⁴ When pressed further on
12 whether OpenAI used YouTube videos, Murati struggled to answer and repeatedly refused to clarify.
13 *Id.* This evasion was deliberate. When OpenAI COO Brad Lightcap was later given another
14 opportunity to clarify if YouTube videos were used to train Sora, OpenAI again refused to answer,
15 opting for silence shrouded in a heap of words.⁵

16 65. The label “publicly available” does not mean lawfully obtained. The fact that a video
17 can be accessed and viewed by the public through a web browser does not mean it can be accessed,
18 downloaded, extracted, and fed into a commercial AI training pipeline. YouTube’s content is
19 publicly viewable precisely because YouTube streams it under strict access controls that prohibit
20 downloading. YouTube CEO Neal Mohan made this distinction explicit, stating that YouTube’s
21 terms of service “does not allow for things like transcripts or video bits to be downloaded, and that
22 is a clear violation of our terms of service.”⁶ Defendant cannot hide behind the word “public” to

23 ⁴ Joanna Stern, *OpenAI CTO Mira Murati on Sora, Generative Video*, Wall St. J. (Mar. 2024),
24 <https://www.wsj.com/tech/personal-tech/openai-cto-sora-generative-video-interview-b66320bb>
25 (when asked directly whether Sora was trained on YouTube videos, Murati stated only that
26 Defendants used “publicly available” and declined to confirm or deny the use of any specific
27 platform's content).

26 ⁵ Jaron Schneider, *OpenAI Again Refuses to Say if It Used Your Content to Train Sora*, PetaPixel
27 (May 10, 2024), [https://petapixel.com/2024/05/10/openai-again-refuses-to-say-if-it-used-your-](https://petapixel.com/2024/05/10/openai-again-refuses-to-say-if-it-used-your-content-to-train-sora/)
28 [content-to-train-sora/](https://petapixel.com/2024/05/10/openai-again-refuses-to-say-if-it-used-your-content-to-train-sora/)

27 ⁶ Davey Alba & Emily Chang, *YouTube Says OpenAI Training Sora With Its Videos Would Break*
28 *Rules*, Bloomberg (Apr. 4, 2024), [https://www.bloomberg.com/news/articles/2024-04-04/youtube-](https://www.bloomberg.com/news/articles/2024-04-04/youtube-says-openai-training-sora-with-its-videos-would-break-the-rules)
[says-openai-training-sora-with-its-videos-would-break-the-rules](https://www.bloomberg.com/news/articles/2024-04-04/youtube-says-openai-training-sora-with-its-videos-would-break-the-rules) (video available at
https://www.youtube.com/watch?v=FBZ_BeChRg).

1 launder what was in reality an unauthorized mass unauthorized access and extraction of protected
2 content.

3 66. Sora is not a research tool. It is Defendant’s commercialized text-to-video generative
4 AI model, monetized through the Sora app and integrated across Defendant’s broader product
5 ecosystem. Because Sora was commercialized, Defendant had an overwhelming incentive to acquire
6 training data on an unprecedented scale. Rather than negotiate for lawful licenses, Defendant broke
7 through YouTube’s access protections to obtain the underlying video files necessary to fuel Sora
8 and, by extension, Defendant’s broader product lines.

9 67. Although Defendant has refused to identify their training data sources with
10 specificity, their own Sora System Card provides a critical admission. As Defendant stated: “Sora
11 was trained on diverse datasets, including a mix of publicly available data, proprietary data accessed
12 through partnerships, and custom datasets developed in-house. These consist of: Select publicly
13 available data, mostly collected from industry-standard machine learning datasets and web crawls.”⁷
14 Defendant thus admits that Sora was trained on what the AI industry recognizes as “industry-
15 standard machine learning datasets.”

16 68. That admission is significant because the field of video-based AI model training has
17 a well-established set of industry-standard datasets. On information and belief, among the industry-
18 standard machine learning datasets on which Defendant trained Sora were HD-VG-130M, HD-
19 VILA-100M, Panda-70M, and HowTo100M. Each of these datasets consists entirely of references
20 and location identifiers pointing to millions of YouTube videos. None contains the underlying
21 audiovisual files themselves. Each therefore required Defendant to access and download the
22 referenced videos directly from YouTube. Meaning that using these datasets for commercial AI
23 training necessarily required the mass circumvention of YouTube's TPMs.

24 69. HD-VG-130M, HD-VILA-100M, Panda-70M and HowTo100M are not obscure or
25 experimental datasets. They are among the most widely cited and broadly adopted video datasets in
26 the field of AI model development. HowTo100M, published by Miech et al., introduced a large-

27 ⁷ OpenAI, *Sora System Card*, OpenAI (Dec. 2024), <https://openai.com/index/sora-system-card/>
28 (disclosing that Sora was trained on “a mix of publicly available data” consisting of “select publicly
available data, mostly collected from industry-standard machine learning datasets and web crawls”
without identifying the specific sources of publicly available data).

1 scale dataset of 136 million video clips sourced from 1.22 million narrated instructional YouTube
2 videos and demonstrated that text-video embeddings trained on this data produced state-of-the-art
3 results across multiple AI benchmarks.⁸ HD-VILA-100M, published by researchers at Microsoft
4 Research Asia, described itself as the largest and most diversified high-resolution video-language
5 dataset available, comprising 100 million clip-sentence pairs drawn from 3.3 million YouTube
6 videos across 15 popular content categories.⁹ HD-VG-130M followed as a large-scale dataset
7 compiled specifically for text-to-video generation research and released publicly on GitHub.¹⁰
8 Panda-70M is a collection of 3.8 million videos from YouTube split into approximately 70.7 million
9 clips and paired with text captions. It was compiled by Snap and released in 2024.¹¹ All four datasets
10 are routinely cited across the AI research literature as foundational resources for training video
11 generation models, and all three are built substantially from YouTube content.

12 70. The centrality of these four datasets to the field of video AI is not a matter of dispute.
13 Independent researchers, academic institutions, and AI companies around the world have cited and
14 built upon HowTo100M, HD-VILA-100M, Panda-70M and HD-VG-130M as standard reference
15 points for years. They are, in the plain meaning of the term Defendant used in their Sora System
16 Card, “industry-standard machine learning datasets.” When Defendant admitted that Sora was
17 trained on industry-standard machine learning datasets collected from the web, they admitted that
18 Sora was trained on datasets of precisely this kind. The logical inference is direct and supported by
19 the evidence: Defendant used these datasets, built from millions of YouTube videos belonging to
20 Plaintiffs and Class Members, to train Sora commercially, in violation of the express restrictions
21 governing each dataset and in circumvention of YouTube's TPMs.

22 71. At the time OpenAI developed and trained Sora, the use of Panda-70M and HD-

23 ⁸ Antoine Miech et al., *HowTo100M: Learning a Text-Video Embedding by Watching Hundred*
24 *Million Narrated Video Clips*, arXiv:1906.03327 (June 7, 2019), <https://arxiv.org/abs/1906.03327>
(introducing a large-scale dataset of 136 million video clips sourced from 1.22 million narrated
25 instructional YouTube videos spanning over 23,000 visual tasks).

26 ⁹ Hongwei Xue et al., *Advancing High-Resolution Video-Language Representation with Large-*
27 *Scale Video Transcriptions*, arXiv:2111.10337 (Nov. 19, 2021), <https://arxiv.org/abs/2111.10337>.

28 ¹⁰ Wenjing Wang et al., *VideoFactory: Swap Attention in Spatiotemporal Diffusions for Text-to-*
Video Generation, arXiv:2305.10874 (May 2023), <https://arxiv.org/abs/2305.10874> (introducing
the HD-VG-130M dataset, consisting of 130 million text-video pairs collected from YouTube in
high-definition, widescreen, and watermark-free formats for text-to-video generation research).

¹¹ Tsai-Shien Chen et al., *Panda-70M: Captioning 70M Videos with Multiple Cross-Modality*
Teachers, arXiv:2402.19479 (Feb. 29, 2024), <https://arxiv.org/abs/2402.19479>.

1 VG-130M represented the industry standard for large-scale text-to-video model training. This is
2 demonstrated by the independent conduct of multiple research teams who, lacking access to
3 OpenAI's proprietary pipeline, built Sora-equivalent models using precisely these datasets. The
4 PKU-YuanGroup's Open-Sora-Plan, initiated by Peking University researchers explicitly to
5 reproduce Sora, trained its models primarily on Panda-70M across every version of the project,
6 treating it as the foundational video corpus for achieving Sora-class output quality.¹²

7 72. In another example of an independent research group attempting to recreate Sora,
8 the “Open-Sora” project likewise trained on Panda-70M and HD-VG-130M, identifying these
9 datasets as core components of the open-source training corpus necessary to reproduce Sora’s
10 capabilities.¹³ The convergent, independent reliance on these specific datasets by multiple
11 sophisticated research teams confirms that they were not merely one option among many. These are
12 datasets a developer would use to build a model like Sora. A company with OpenAI's resources
13 building a model of Sora's capability would have used, at minimum, the same datasets the open-
14 source community identified as necessary, and in all likelihood far more of them.

15 73. YouTube allows users to stream content but prohibits, through contractual terms and
16 TPMs, accessing the underlying video file of the content it streams. Defendant obtained
17 unauthorized access to audiovisual content by unlawfully accessing and stream ripping, converting
18 streaming content into permanent, locally stored files, and by circumventing the TPMs specifically
19 designed to prevent such access.

20 74. Each of the four datasets functions as a map or index file identifying specific
21 YouTube videos and clips by URL, video identifier, and timestamp. A single YouTube video may
22 be divided into numerous clips, each treated as a separate training sample. Extracting any clip
23 requires independently accessing the source video on YouTube and isolating the designated
24 segment, a process that constitutes a separate act of circumvention for each clip retrieved.

25 75. Defendant obtained unauthorized access to audiovisual content to train its AI models
26 by unlawfully “stream ripping” them—*i.e.*, converting the audiovisual streaming content into

27 ¹² Bin Lin et al., *Open-Sora Plan: Open-Source Large Video Generation Model*,
arXiv:2412.00131 (Nov. 28, 2024), <https://arxiv.org/abs/2412.00131>.

28 ¹³ Zangwei Zheng et al., *Open-Sora: Democratizing Efficient Video Production for All*,
arXiv:2412.20404 (Dec. 29, 2024), <https://arxiv.org/abs/2412.20404>.

1 permanent, locally-stored files—from YouTube and circumventing TPMs specifically designed to
2 prevent access for such use.

3 76. The HD-VG-130M dataset contains references to more than 130 million video clips
4 derived from 1,549,408 YouTube videos. The dataset was created by partitioning approximately
5 1,549,408 YouTube videos into short high-definition segments intended for machine learning
6 research. Each listed clip corresponds to a specific portion of a source video that must be
7 independently retrieved in order to obtain usable audiovisual data.

8 77. The compilers of HD-VG-130M selected videos based on automated criteria,
9 including aesthetic scoring, and excluded videos displaying creator names or visible watermarks.
10 The dataset was released for research purposes and does not itself distribute the underlying videos
11 or grant any license to use those works commercially.

12 78. Defendant used HD-VG-130M to retrieve every clip and incorporate those files into
13 its commercial training pipeline. Because each clip corresponds to a separate act of circumvention,
14 the ingestion of HD-VG-130M required roughly 130 million unauthorized downloads, each
15 constituting a separate circumvention event.

16 79. Microsoft created a document identifying the content of the HD-VG-130M dataset.
17 Microsoft then released this list into the broader AI ecosystem. Not only did Microsoft assemble the
18 HD-VG-130M dataset internally, but it also posted the dataset publicly on GitHub—a code-sharing
19 platform owned by Microsoft—at <https://github.com/daoshee/HD-VG-130M>.

20 80. The GitHub repository includes a license agreement that restricts the dataset to
21 “academic use only,” acknowledging that any use for commercial purposes would be unauthorized.

22 81. In addition, the GitHub license agreement prohibits distribution, reproduction,
23 modification, or exploitation of the dataset content without the permission of the copyright owner,
24 thus recognizing that copyright ownership resides with the original creators of the source material,
25 including YouTube content creators such as Plaintiffs and the Class Members.

26 82. By including these limitations on its GitHub repository, Microsoft implicitly
27 acknowledged that the HD-VG-130M dataset is comprised of protected, copyrighted works obtained
28 without license, permission, or authorization.

1 83. On information and belief, Defendant obtained the dataset that was earmarked for
2 academic use only, including the HD-VG-130M dataset. Defendant was aware that the datasets were
3 for research purposes only but used the datasets for commercial development (*i.e.*, “training”) of its
4 generative AI model.

5 84. Further, Defendant used the HD-VG-130M dataset to scrape audiovisual content
6 from YouTube without the consent of YouTube or the content creators.

7 85. The HD-VILA-100M dataset is comprised of references to roughly 100 million clips
8 that were derived from 3,098,462 YouTube videos. Each clip represents a distinct segment of the
9 original work. The dataset was compiled by Microsoft Research Asia and published in 2021 for
10 research involving video-based AI models. It lists pointers to YouTube videos and includes subtitle
11 information extracted from YouTube’s closed captioning system.

12 86. The usage license on the dataset states it was created “solely for Computational Use
13 for non-commercial research. This restriction means that you may engage in non-commercial
14 research activities (including non-commercial research undertaken by or funded via a commercial
15 entity), but you may not use the Data or any Results in any commercial offering, including as part
16 of a product or service (or to improve any product or service) you use or provide to others.”

17 87. HD-VILA-100M does not contain the video files themselves. Any user of the dataset
18 must download each referenced clip directly from YouTube. Defendant did so at scale, which
19 triggered roughly 100 million separate acts of unauthorized access and copying.

20 88. By including these limitations on its usage license, Microsoft implicitly
21 acknowledged that the HD-VILA-100M dataset is comprised of protected, copyrighted works
22 obtained without license, permission, or authorization.

23 89. Defendant obtained the dataset that was earmarked for non-commercial use only,
24 including the HD-VILA-100M dataset. Defendant was aware that the datasets were for research
25 purposes only but used the datasets for commercial development (*i.e.*, “training”) of its generative
26 AI model. Further, Defendant used the HD-VILA-100M dataset to scrape audiovisual content from
27 YouTube without the consent of YouTube or the creators of the audiovisual content.

28 90. The HowTo100M dataset consists of references to more than 100 million clips that

1 were derived from 1,238,911 YouTube videos. The clips depict 23,611 different tasks related to
2 cooking, handcrafting, personal care, gardening, and other activities. The task labels and
3 descriptions were generated using the taxonomy developed by WikiHow. The dataset is hosted by
4 the European research institutions that created it. Like the others, it contains only identifiers and
5 metadata, not any audiovisual material. To use HowTo100M in training, a company must download
6 each clip directly from YouTube. Defendant carried out these downloads as part of its Sora training
7 pipeline, resulting in more than 100 million additional unauthorized copying events.

8 91. Although the HowTo100M dataset is posted online with what its creators label as a
9 “commercial license,” that license has no effect on the legal rights associated with the underlying
10 YouTube videos. The dataset creators do not own the videos and therefore cannot grant any license
11 that authorizes the reproduction, downloading, or commercial use of those works. The dataset
12 license applies only to the dataset itself, meaning the text files, identifiers, and metadata created by
13 the researchers. It does not and cannot confer any right to copy, download, or redistribute the
14 millions of YouTube videos identified in the dataset. Each of those videos remains protected by
15 YouTube’s Terms of Service and by technological access restrictions. As a result, every act of
16 downloading a HowTo100M clip from YouTube is a separate unauthorized accessing and copying
17 event and a separate circumvention violation, regardless of any license attached to the dataset file.

18 92. Panda-70M is a derivative dataset created by Snap that used the HD-VILA-100M
19 dataset as the foundation for its own refined dataset. Panda-70M is a collection of 3.8 million videos
20 from YouTube split into approximately 70.7 million clips and paired with text captions. It was
21 compiled by Snap and released in 2024. Developers used AI to create a new set of captions
22 describing what is pictured in each clip.

23 93. The HDVILA-100M, HD-VG-130M, Panda-70M, and HowTo100M datasets
24 consist of location identifiers that point to millions of YouTube videos or clips. None of them
25 contains the underlying audiovisual files. To use them in training, a company must access, retrieve
26 and download every referenced video directly from YouTube. Defendant used these datasets to
27 initiate millions of individual downloads of protected YouTube content, all without authorization or
28 compensation, all in violation of YouTube’s TPMs, and all for the commercial purpose of building

1 its foundational video model.

2 **D. Defendant’s Automated Scraping Required Circumvention of YouTube’s**
3 **Technological Protection Measures.**

4 94. Bulk extraction of YouTube videos at the scale alleged here cannot occur without
5 circumventing YouTube’s layered technological protection measures (“TPMs”), each of which
6 independently controls access to the underlying audiovisual files. Illicit tools and services are
7 specifically designed to defeat these protections, enabling automated systems to retrieve permanent
8 copies of content that YouTube makes available only for streaming. This process is commonly
9 known as “scraping” or “stream ripping.” Under 17 U.S.C. § 1201(a), circumvention of any single
10 technological measure that effectively controls access to a copyrighted work is sufficient to establish
11 liability; YouTube’s protections are layered and overlapping, and the circumvention of one does not
12 negate the effectiveness or legal significance of the others.

13 95. “Scraping” or “stream ripping” content involves the automated accessing and
14 downloading of audio and video files directly from a website’s servers rather than viewing the
15 content through the provider’s authorized playback environment.

16 96. Upon information and belief, in order to acquire “high-quality text-to-video
17 generation,” Defendant, directly and through agents, contractors, and affiliates, intentionally
18 accessed large volumes of YouTube videos by scraping and/or using tools and workflows that
19 bypass or evade YouTube’s TPMs and usage restrictions, and then reproduced those videos to
20 assemble training corpora for Defendant’s AI models and services.

21 97. When Defendant scraped audio and video files from YouTube, Defendant did not
22 simply access and download those files onto the YouTube app for offline streaming, as envisioned
23 by YouTube’s Premium plan. Instead, Defendant improperly accessed the underlying audio and
24 video files and downloaded those files into Defendant’s own system, where Defendant had access
25 over the files and where Defendant could store them indefinitely. These actions circumvented
26 YouTube’s TPMs, violated YouTube’s Terms of Service, and are inconsistent with the access
27 provided by YouTube’s Premium plan.

28 98. Defendant also obtained its own separate audiovisual datasets directly through its

1 own independent scraping of YouTube's video sharing platform.

2 99. To retrieve audiovisual files at scale, Defendant was required to defeat TPM (1), the
3 rolling cipher that protects the true media file URL. On information and belief, Defendant used one
4 or more descrambling tools designed to defeat YouTube's proprietary signature-transformation
5 logic, allowing automated systems to generate valid file requests outside the authorized player
6 environment and thereby obtain the underlying media files directly. Such tools include, for example,
7 the open-source program yt-dlp, which is specifically engineered to descramble, replicate, or extract
8 YouTube's proprietary signature-transformation logic and circumvent the access controls YouTube
9 uses to prevent unauthorized downloading.

10 100. A descrambling tool such as yt-dlp works by replicating or extracting the signature-
11 transformation process that YouTube uses to protect its media URLs, allowing a user to bypass the
12 authorized streaming environment entirely and retrieve the underlying video and audio files directly
13 from YouTube's servers. Tools of this kind can be used to automatically merge separate audio and
14 video streams and to download entire playlists at scale. On information and belief, Defendant used
15 yt-dlp or a substantially similar descrambling tool for precisely this purpose.

16 101. To deploy a descrambling tool such as yt-dlp, a user must first obtain and configure
17 the program, along with ancillary software necessary to merge video and audio streams into
18 complete audiovisual files. Once configured, the tool can be directed at individual videos or entire
19 playlists by supplying the Uniform Resource Locator ("URL") for each target video. The result is a
20 complete audiovisual file retrieved outside of and in circumvention of YouTube's authorized
21 playback environment.

22 102. In simpler terms, a descrambling tool such as yt-dlp acts as a bootleg key for the
23 lock imposed by YouTube's rolling cipher. On information and belief, Defendant used yt-dlp or a
24 substantially similar tool to improperly access and download audiovisual files from YouTube and
25 merge them into complete audiovisual packages for ingestion into their generative AI training
26 pipeline.

27 103. Defendant also defeated TPM (2), YouTube's IP-based monitoring and blocking
28 system, which restricts high-volume accessing and downloading from a single source. On

1 information and belief, Defendant implemented an IP-rotation scheme so that when YouTube
2 detected automated activity and blocked one address, accessing and downloading could immediately
3 resume from another. This allowed continuous access to audiovisual files despite YouTube's efforts
4 to detect and terminate abusive activity.

5 104. YouTube uses automated programs to monitor activity from IP addresses, detect
6 high-volume accessing and downloading behavior, and block those IP addresses from further access
7 to audiovisual content on the platform. An IP-rotation scheme is specifically designed to defeat this
8 monitoring system by cycling through different IP addresses to avoid detection and blocking.

9 105. Executing an IP-rotation scheme at the scale necessary to download the volume of
10 videos required to train Sora necessarily required the use of virtual machines or substantially similar
11 infrastructure. A single physical machine with a fixed IP address would have been detected and
12 blocked almost immediately given the volume of requests involved. On information and belief,
13 Defendant deployed virtual machines, cloud computing infrastructure, or substantially similar
14 technology to rotate IP addresses at scale, ensuring that when YouTube detected and blocked one
15 address, Defendant's downloading operation could immediately resume from another.

16 106. On information and belief, by repeatedly cycling through IP addresses in this
17 manner, Defendant ensured that their automated downloading operation could continue
18 uninterrupted even after YouTube detected and blocked individual addresses. The sheer volume of
19 videos necessary to train Sora makes any alternative explanation implausible. The full details of the
20 specific infrastructure Defendant used to execute this scheme will be the subject of discovery.

21 107. Processes such as these allowed Defendant to bypass YouTube's player page and
22 avoid YouTube's monitoring systems in order to access and scrape content from YouTube.

23 108. Defendant's conduct necessarily circumvented TPM (3), the use of session-bound,
24 short-lived media URLs. Because these URLs expire and cannot be reused, automated tools must
25 repeatedly obtain fresh authorization parameters and regenerate valid links in order to continue
26 downloading. By programmatically renewing access credentials outside ordinary playback sessions,
27 Defendant maintained uninterrupted access to files that would otherwise become unavailable once
28 the original authorization lapsed.

1 109. At scale, Defendant’s automated activity would also trigger TPM (4), CAPTCHA
2 challenges intended to distinguish human users from automated systems. To continue retrieving
3 videos without interruption, Defendant’s infrastructure distributed requests across numerous
4 machines and operated in a manner designed to avoid or neutralize such verification barriers, thereby
5 obtaining audiovisual data without the human interaction required for continued access.

6 110. Defendant further circumvented TPM (5), YouTube’s proof-of-origin token system,
7 which requires requests for video segments to present credentials demonstrating that they originate
8 from an authorized playback environment. Automated downloading tools operate outside that
9 environment and therefore must extract, replicate, or reuse the necessary request parameters in order
10 to retrieve video data directly, allowing Defendant to access and obtain files that would otherwise
11 be delivered only to approved clients.

12 111. Defendant did not obtain the consent of YouTube, Plaintiffs or other Class Members
13 to conduct its scraping activities.

14 112. On information and belief, and based on the sheer volume of video data included,
15 the datasets of videos that included YouTube videos were scraped without permission from
16 YouTube and in violation of YouTube’s Terms and Services.

17 113. The automated system used by Defendant to access video content from YouTube
18 was intentionally designed and utilized to circumvent YouTube’s TPMs.

19 114. To be clear, Defendant’s generative AI models are not “watching” YouTube in order
20 to be trained. Rather, data improperly accessed and downloaded from YouTube is being uploaded
21 into Defendant’s generative AI models to develop and improve Defendant’s products.

22 115. Defendant knew or should have known its conduct contravened YouTube’s rules
23 and the TPMs that enforce them.¹⁴

24 116. Defendant’s actions constitute a breach of the Copyright Act’s anti-circumvention
25 provisions, which state, among other things, that “[n]o person shall circumvent a technological
26 measure that effectively controls access to a work protected under this title.” 17 U.S.C. §

27 ¹⁴ YouTube’s Terms of Service specifically state that users “may access and use the Service as
28 made available to you, *as long as you comply with this Agreement* and applicable law.” “Terms of
Service,” YouTube, <https://www.youtube.com/t/terms> (last viewed March 12, 2026).

1 1201(a)(1)(a).

2 **E. Defendant's Actions Caused Harm to Plaintiffs and Class Members**

3 117. At no time did Defendant seek or obtain Plaintiffs' or Class Members' authorization
4 to access and use the videos Defendant used to train its generative AI models.

5 118. At no time did Defendant seek or obtain Plaintiffs' or Class Members' authorization,
6 nor did Defendant compensate Plaintiffs or the Class Members for the access, copying, ingestion,
7 and use of their YouTube videos in AI training and related workflows.

8 119. Plaintiffs and Class Members used YouTube as their platform of choice to upload
9 their video content in substantial part due to YouTube's Terms of Service that prohibit the type of
10 scraping, unauthorized accessing and downloading, bulk extraction, and other forms of data mining
11 of audiovisual content utilized by Defendant to obtain Plaintiffs' video content.

12 120. The harm Defendant is causing goes beyond the immediate economic consequences
13 of unauthorized and uncompensated extraction and use of Plaintiffs' and Class Members' works. It
14 degrades the rights of creators to control their works, determine whether future uses of their work
15 align with their values, and decide the products or services with which they wish to be associated.

16 **CLASS ACTION ALLEGATIONS**

17 121. Plaintiffs bring this action on behalf of themselves and on behalf of all other persons
18 similarly situated.

19 122. Plaintiffs propose the following Class definition, subject to amendment as
20 appropriate:

21 **All individuals and entities in the United States whose videos were posted on**
22 **YouTube and were accessed by Defendant.**

23 123. Excluded from the Class are Defendant's officers and directors, and any entity in
24 which Defendant has a controlling interest; and the affiliates, legal representatives, attorneys,
25 successors, heirs, and assigns of Defendant. Excluded also from the Class are Members of the
26 judiciary to whom this case is assigned, their families and Members of their staff.

27 124. Plaintiffs reserve the right to amend or modify the class definition with greater
28 specificity or division after having an opportunity to conduct discovery. The proposed Class meets

1 the criteria for certification under Rule 23 of the Federal Rules of Civil Procedure.

2 125. Plaintiffs also reserve the right to seek certification of subclasses, or alternative
3 classes as shall become necessary in the course of litigation.

4 126. Numerosity. The Members of the Class are so numerous that joinder of all of them
5 is impracticable. The exact number of Class Members is unknown to Plaintiffs now, but Plaintiffs
6 estimates that there are thousands of Class Members.

7 127. Commonality. There are questions of law and fact common to the Class, which
8 predominate over any questions affecting only individual Class Members. These common questions
9 of law and fact include, without limitation:

10 a. Whether Defendant used an automated tool to download or otherwise obtain
11 unauthorized access to the Plaintiffs' and Class Members' copyrighted content on YouTube;

12 b. Whether Defendant's downloading of the video content was intended to bypass the
13 TPMs put in place by YouTube at the authorization of Plaintiffs and Class Members;

14 c. Whether Defendant used Plaintiffs' and Class Members' original content for
15 purposes of training its generative AI models;

16 d. Whether Plaintiffs and Class Members suffered damages from Defendant's
17 misconduct;

18 e. Whether Plaintiffs and Class Members are entitled to actual damages, statutory
19 damages and/or injunctive relief.

20 128. These issues are common to the Class, and their resolution would advance this
21 matter and the parties' interests therein.

22 129. Typicality. Plaintiffs' claims are typical of those of other Class Members because
23 Plaintiffs' video content, like that of every other Class Member, was made available on YouTube
24 and improperly altered and used by Defendant to train Defendant's generative AI models for
25 commercial purposes. Plaintiffs' claims are typical of those of the other Class Members because,
26 among other things, all Class Members were injured through the common misconduct of Defendant.
27 Plaintiffs are advancing the same claims and legal theories on behalf of themselves and all other
28 Class Members, and no defenses are unique to Plaintiffs. Plaintiffs' claims and those of Class

1 Members arise from the same operative facts and are based on the same legal theories.

2 130. Adequacy of Representation. Plaintiffs will fairly and adequately represent and
3 protect the interests of the Members of the Class. Plaintiffs' Counsel is competent and experienced
4 in litigating class actions.

5 131. Predominance. Defendant has engaged in a common course of conduct toward
6 Plaintiffs and Class Members, in that Plaintiffs' and Class Members' video content was unlawfully
7 downloaded, altered and used by Defendant in the same way. The fact that Sora was a foundational
8 model underscores that Defendant's conduct was uniform across the Class: every video unlawfully
9 accessed from YouTube was ingested into the same core model that underpins Defendant's product
10 ecosystem. This commonality further demonstrates the predominance of shared legal and factual
11 issues among Class Members. The common issues arising from Defendant's conduct affecting Class
12 Members set out above predominate over any individualized issues. Adjudication of these common
13 issues in a single action has important and desirable advantages of judicial economy.

14 132. Superiority. A Class action is superior to other available methods for the fair and
15 efficient adjudication of the controversy. Class treatment of common questions of law and fact is
16 superior to multiple individual actions or piecemeal litigation. Absent a class action, most Class
17 Members would likely find that the cost of litigating their individual claims is prohibitively high and
18 would therefore have no effective remedy. The prosecution of separate actions by individual Class
19 Members would create a risk of inconsistent or varying adjudications with respect to individual
20 Class Members, which would establish incompatible standards of conduct for Defendant. In
21 contrast, the conduct of this action as a class action presents far fewer management difficulties,
22 conserves judicial resources and the parties' resources, and protects the rights of each Class member.

23 133. Defendant has acted on grounds that apply generally to the Class as a whole, so that
24 class certification, injunctive relief, and corresponding declaratory relief are appropriate on a Class-
25 wide basis.

26 134. Ascertainability. Finally, all members of the proposed Class are readily
27 ascertainable. Defendant used datasets that contain the full lists of URLs and video identifiers for
28 every YouTube video incorporated into the training pipeline. Those identifiers allow the parties to

1 determine exactly which videos were used and to match each video to its creator through YouTube's
2 public channel and authorship information. Because the datasets provide a complete map of the
3 videos Defendant downloaded, the identities of the affected creators are identifiable through
4 straightforward reference to the URLs and corresponding channel data.

5 **CLAIM FOR RELIEF**

6 **Violation of the DMCA (Anti-Circumvention) 17 U.S.C. § 1201(a)**

7 135. Plaintiffs repeat, reallege, and incorporates the allegations contained in the previous
8 paragraphs as if fully set forth herein.

9 136. YouTube's Terms of Service prohibit scraping, unauthorized downloading, bulk
10 extraction, or other forms of data mining of audiovisual content. These restrictions operate together
11 with TPMs to prevent unlicensed access to the underlying video files, not simply reproduction of
12 content.

13 137. Plaintiffs and Class Members uploaded their original video content to YouTube's
14 video sharing platform due in substantial part to YouTube's Terms of Service and TPMs prohibiting
15 bulk downloading of creators' videos.

16 138. Content creators, including Plaintiffs and Class Members, expect that their works
17 will not be accessed at the file level and copied at scale without consent. Plaintiffs and Class
18 Members own all rights, title, and interest in and to the claims asserted in this action, including all
19 claims for violation of 17 U.S.C. § 1201.

20 139. Each of Plaintiffs' and Class Members works were uploaded to YouTube. When
21 published, YouTube automatically applied TPMs that control access to and prevent unauthorized
22 copying or downloading of these audiovisual files.

23 140. YouTube's Terms of Service and associated access controls constitute "effective
24 technological measures" within the meaning of 17 U.S.C. §1201.

25 141. Defendant used automated tools for the sole purpose of circumventing YouTube's
26 technological barriers that effectively control access to extracted files never made available to the
27 public. Specifically, Defendant employed measures to access, extract, copy, and download
28 Plaintiffs' and Class Members' content from YouTube without authorization. In doing so, Defendant

1 improperly obtained millions of videos from YouTube’s platform.

2 142. This distinction is critical: viewing a YouTube video through YouTube’s authorized
3 and advertising-supported playback mechanisms does not provide access to the underlying file.
4 Defendant’s circumvention tools broke through that access barrier, triggering liability under §1201.

5 143. Plaintiffs and Class Members have been harmed by Defendant’s violations of 17
6 U.S.C. §1201(a) because Defendant has taken their content without authorization or compensation.
7 Defendant’s circumvention of YouTube’s technological measures has facilitated Defendant’s
8 ongoing and mass-scale copyright infringement through its unauthorized and uncompensated use of
9 the content in its training data.

10 144. Defendant accessed, downloaded, stored and utilized those videos to assemble
11 training corpora for Defendant’s AI models and services.

12 145. By using circumvention software and disabling or bypassing YouTube’s TPMs,
13 Defendant violated 17 U.S.C. § 1201(a), which prohibits the circumvention of a technological
14 measure that effectively controls access to a copyrighted work.

15 146. Each act of circumvention constitutes a separate violation of §1201. Plaintiffs and
16 the Class Members are entitled to statutory damages, injunctive relief, impoundment, and attorneys’
17 fees and costs under 17 U.S.C. §1203.

18 147. Each of Defendant’s acts of infringement is a willful violation as Defendant
19 specifically utilized automated tools designed to evade YouTube’s TPMs.

20 **PRAYER FOR RELIEF**

21 WHEREFORE, Plaintiffs respectfully request a judgment in their favor and against
22 Defendant as follows:

23 A. For a declaration that Defendant has willfully circumvented the copyright
24 protection systems of YouTube intended to protect Plaintiffs’ and the Class Members’s
25 audiovisual content.

26 B. For statutory damages (up to the maximum allowed by law per violation),
27 injunctive relief, and attorneys’ fees and costs under 17 U.S.C. §1203;

28 C. For such equitable relief under Title 17, Title 28, and/or the Court’s inherent

1 authority as is necessary to prevent or restrain infringement of Plaintiffs’ and the Class
2 Members’ copyright-protected content, including a preliminary and permanent injunction
3 requiring that Defendant and its officers, agents, servants, employees, attorneys, directors,
4 successors, assigns, licensees, and all others in active concert or participation with any of them,
5 cease infringing, or causing, aiding, enabling, facilitating, encouraging, promoting, inducing,
6 or materially contributing to or participating in the infringement of any of Plaintiff’s or the
7 Class Members’ exclusive rights under federal law, including without limitation in the content
8 identified in Exhibits A, B, and C;

9 D. For an award of pre-judgment and post-judgment interest, to the fullest extent
10 available, on any monetary award made part of the judgment against Defendant; and

11 E. For such other and further relief as the Court may deem just and proper.

12 **JURY DEMAND**

13 Plaintiffs demand a trial by jury on all claims for which trial by jury is proper.

14 Dated: April 3, 2026

HEAH BAR-NISSIM LLP

15
16 By /s/ ROM BAR-NISSIM
ROM BAR-NISSIM

**ELLZEY KHERKHER SANFORD
MONTGOMERY LLP**

JARRET LEE ELLZEY
TOM KHERKHER
LEIGH S. MONTGOMERY

21 Attorneys for Plaintiffs