

AMENDED IN SENATE AUGUST 23, 2024
AMENDED IN SENATE AUGUST 22, 2024
AMENDED IN SENATE JUNE 24, 2024
AMENDED IN SENATE JUNE 10, 2024
AMENDED IN ASSEMBLY APRIL 18, 2024
AMENDED IN ASSEMBLY MARCH 21, 2024
CALIFORNIA LEGISLATURE—2023–24 REGULAR SESSION

ASSEMBLY BILL

No. 3211

Introduced by Assembly Member Wicks

February 16, 2024

An act to add Chapter 41 (commencing with Section 22949.90) to Division 8 of the Business and Professions Code, relating to artificial intelligence.

LEGISLATIVE COUNSEL'S DIGEST

AB 3211, as amended, Wicks. California Digital Content Provenance Standards.

Existing law requires the Secretary of Government Operations to develop a coordinated plan to, among other things, investigate the feasibility of, and obstacles to, developing standards and technologies for state departments to determine digital content provenance. For the purpose of informing that coordinated plan, existing law requires the secretary to evaluate, among other things, the impact of the proliferation of deepfakes, as defined.

This bill, the California Digital Content Provenance Standards, would require a generative artificial intelligence (AI) provider, as provided,

to, among other things, apply provenance data to synthetic content produced or significantly modified by a generative AI system that the provider makes available, as those terms are defined, and to conduct adversarial testing exercises, as prescribed. The bill would prohibit, among other things, providers and distributors of software and online services from making available a system, application, tool, or service that is designed for the primary purpose of removing provenance data from synthetic content, as provided.

This bill would require a newly manufactured recording device sold, offered for sale, or distributed in California to offer users the option to apply difficult to remove provenance data to nonsynthetic content produced by that device and would require the application of that provenance data to be compatible with state-of-the-art adopted and relevant industry standards. If technically feasible and secure, the bill would require a recording device manufacturer to offer a software or firmware update enabling a user of a recording device manufactured before July 1, 2026, and purchased in California to apply difficult to remove provenance data to the nonsynthetic content created by the device and decode any provenance data attached to nonsynthetic content created by the device.

This bill would require a large online platform, as defined, capable of disseminating specified content to use labels to disclose, as specified, any machine-readable provenance data detected in synthetic content that is distributed on its platform. If content uploaded to or distributed on a large online platform by a user does not contain specified provenance data or if the content's provenance data cannot be interpreted or detected, the bill would require a large online platform to label the content as having unknown provenance. The bill would require a large online platform to use a visual disclosure that contains specified information, including the copyright holder or licensor information, when labeling and disclosing provenance data of sound recordings and music videos.

Beginning July 1, 2026, and annually thereafter, this bill would require a large online platform to produce a transparency report that identifies moderation of deceptive synthetic content on their platform and would authorize that report to include, among other things, instances where synthetic or potentially deceptive content was identified and removed by the platform, as applicable.

This bill would authorize the Department of Technology (department) to assess specified administrative penalties for prescribed violations of

the bill's provisions, including an administrative penalty of up to ~~\$500,000~~ \$100,000 for each violation that is intentional or is the result of grossly negligent conduct, to be deposited in the Digital Content Provenance Administrative Fund which the bill would establish in the State Treasury. The bill would, upon appropriation by the Legislature for this express purpose, authorize the expenditure of moneys in the fund by the department to administer these provisions.

This bill would make its provisions operative on July 1, 2026.

Vote: majority. Appropriation: no. Fiscal committee: yes.

State-mandated local program: no.

The people of the State of California do enact as follows:

1 SECTION 1. The Legislature finds and declares all of the
2 following:

3 (a) Generative artificial intelligence (GenAI) technologies are
4 increasingly able to synthesize images, audio, and video content
5 in ways that are harmful to society.

6 (b) In order to reduce the severity of the harms caused by GenAI,
7 it is important for photorealistic synthetic content to be clearly
8 disclosed and labeled.

9 (c) Failing to appropriately label synthetic content created by
10 GenAI technologies can skew election results, enable defamation,
11 and erode trust in the online information ecosystem.

12 (d) The Legislature should act to adopt standards pertaining to
13 the clear disclosure and labeling of synthetic content, in order to
14 alleviate harms caused by the misuse of GenAI technologies.

15 (e) The Legislature should push for the creation of tools that
16 allow Californians to assess the provenance of content distributed
17 online and the ways in which content has been significantly altered
18 or completely synthesized by GenAI.

19 (f) The Legislature should require online platforms to label
20 synthetic content produced by GenAI.

21 (g) Through these actions, the Legislature can help to ensure
22 that Californians remain safe and informed.

23 SEC. 2. Chapter 41 (commencing with Section 22949.90) is
24 added to Division 8 of the Business and Professions Code, to read:

CHAPTER 41. CALIFORNIA DIGITAL CONTENT PROVENANCE
STANDARDS

22949.90. For purposes of this chapter, the following definitions apply:

(a) “Adversarial testing” means a structured testing effort to find flaws and vulnerabilities in a generative AI system’s ability to attach robust provenance data to synthetic content created by the system and access potential risks associated with misuse of the generative AI system to attach false provenance data to digital content generated outside of the generative AI system.

(b) “Artificial intelligence” or “AI” means an engineered or machine-based system that varies in its level of autonomy and that can, for explicit or implicit objectives, infer from the input it receives how to generate outputs that can influence physical or virtual environments.

(c) “Digital fingerprint” means a unique value that can be used to identify identical or similar digital content.

(d) “Digital signature” means a cryptography-based method that identifies the user or entity that attests to the information provided in the signed section.

(e) “Generative AI hosting platform” means an online repository or other internet website that makes a generative AI system available for use by a California resident, regardless of whether the terms of that use include compensation.

(f) “Generative AI provider” or “GenAI provider” means an organization or individual that creates, codes, substantially modifies, or otherwise produces a generative AI system that is made publicly available for use by a California resident, regardless of whether the terms of that use include compensation.

(g) “Generative AI system” or “GenAI system” means an artificial intelligence system that can generate derived synthetic content, including images, videos, and audio, and that emulates the structure and characteristics of the system’s training data.

(h) “Large online platform” means a public-facing social media platform, as defined in Section 22675, video-sharing platform, messaging platform, advertising network, or standalone search engine that displays content to viewers who are not the creator or collaborator and had at least 2,000,000 unique monthly California users during the preceding 12 months.

1 (i) “Metadata” means structural or descriptive information about
2 data.

3 (j) “Nonsynthetic content” means images, videos, or audio
4 captured in the physical world by natural persons using a recording
5 device and is without any modifications or with only minor
6 modifications that do not lead to significant changes to the
7 perceived contents or meaning of the content. Minor modifications
8 include, but are not limited to, changes to brightness or contrast
9 of images and removal of background noise in audio.

10 (k) “Provenance data” means data that records the origin or
11 history of digital content and is communicated using state-of-the-art
12 techniques based on widely adopted and relevant industry
13 standards. “Provenance data” may be communicated using digital
14 fingerprinting to associate metadata with digital content, attaching
15 metadata to digital content, including through the use of a digital
16 signature, or embedding of watermarks in digital content.

17 (l) “Provenance detection tool” means a software tool or online
18 service that can read or interpret a watermark, metadata, or digital
19 signature, and output the associated provenance data.

20 (m) “Synthetic content” means ~~information, including~~ images,
21 videos, and audio, that has been produced or significantly modified
22 by a generative AI system.

23 (n) “Watermark” means information covertly embedded into
24 digital content, including image, audio, and video, for the purpose
25 of communicating the provenance, history of modification, or
26 history of conveyance.

27 22949.90.1. (a) A generative AI provider whose GenAI system
28 is capable of producing digital content that would falsely appear
29 to a reasonable person to depict real-life persons, objects, places,
30 entities, or events shall do all of the following:

31 (1) (A) Apply provenance data, either directly or through the
32 use of third-party technology, to synthetic content produced or
33 significantly modified by a generative AI system that the GenAI
34 provider makes available. The GenAI provider shall make the
35 provenance data difficult to ~~remove~~, *remove or disassociate*, taking
36 into account the accuracy of the provenance data, the quality of
37 the content produced or significantly modified by the generative
38 AI system, and widely accepted industry standards on ~~the~~
39 ~~robustness of~~ provenance data.

(B) The application of provenance data to synthetic content, as required by subparagraph (A), shall, at minimum, be difficult to ~~remove~~; *remove or disassociate*, identify the digital content as synthetic, and communicate the following provenance data in order of priority, with clause (i) being the most important, and clause (iv) being the least important:

- (i) The synthetic nature of the content.
- (ii) The name of the generative AI provider.
- (iii) If feasible for the provenance technique used, the time and date the provenance data was applied.

(iv) If applicable and feasible for the provenance technique used, the specific portions of the content that are synthetic.

(2) (A) A generative AI provider shall create and make available to the public a provenance detection tool or permit users to use a provenance detection tool provided by a third party. The provenance detection tool shall be based on broadly adopted industry standards and, if technically feasible, meet the following criteria:

(i) The tool allows a user to assess whether digital content was created or altered by a generative AI system.

(ii) The tool allows a user to determine how digital content was created or altered by a generative AI system.

(iii) The tool outputs any provenance data that is detected in the content.

(iv) The tool is publicly accessible through the generative AI provider's or the third-party's internet website, its mobile application, or an application programming interface, as applicable.

(v) The tool allows a user to upload content or provide a uniform resource locator (URL) linking to online content.

(B) A generative AI provider or third party shall put in place a process to collect user feedback related to the efficacy of the provenance detection tool described in subparagraph (A) and incorporate any feedback into any attempt to improve the efficacy of the tool.

(C) A generative AI provider that creates or makes available a provenance detection tool pursuant to subparagraph (A) may limit access to the decoder to ensure the robustness and security of their provenance data techniques.

1 (3) (A) Conduct adversarial testing exercises following relevant
2 guidelines from the National Institute of Standards and Technology.
3 The adversarial testing exercises shall assess both of the following:

- 4 (i) The robustness of provenance data methods.
5 (ii) Whether the generative AI provider's GenAI systems can
6 be used to add false provenance data to content generated outside
7 of the system.

8 (B) Adversarial testing exercises required by this paragraph
9 shall be conducted before the general audience release of any new
10 tool or method used to apply provenance data to synthetic content
11 produced or significantly modified by a generative AI system that
12 the GenAI provider makes available.

13 (C) In the event that a generative AI provider utilizes a
14 third-party tool or method to apply provenance data, the generative
15 AI provider may rely on the testing conducted by the provider of
16 the third-party tool or method pursuant to paragraph (2).

17 (D) A generative AI provider shall submit full reports of its
18 adversarial testing exercises to the Department of Technology
19 within 90 days of conducting an adversarial testing exercise
20 pursuant to this paragraph. The report shall address any material,
21 systemic failures in a generative AI system related to the erroneous
22 or malicious inclusion or removal of provenance data.

23 (E) (i) Upon the request of an accredited academic institution,
24 a generative AI provider shall make available a summary or report
25 of its adversarial testing exercises.

26 (ii) *The provider may deny a request if providing a summary*
27 *or report to the relevant institution would undermine the robustness*
28 *or security of its provenance data techniques.*

29 (F) This paragraph does not require the disclosure of trade
30 secrets, as defined in Section 3426.1 of the Civil Code.

31 (b) Providers and distributors of software and online services
32 shall not make available a system, application, tool, or service that
33 is designed for the primary purpose of removing provenance data
34 from synthetic content in a manner that would be reasonably likely
35 to deceive a consumer of the origin or history of the content.

36 (c) Generative AI hosting platforms shall not make available a
37 generative AI system that does not allow a GenAI provider, to the
38 greatest extent possible and either directly providing functionality
39 or making available the technology of a third-party vendor, to
40 apply provenance data to content created or substantially modified

1 by the system in a manner consistent with specifications set forth
2 in paragraph (1) of subdivision (a).

3 22949.90.2. (a) (1) A newly manufactured recording device
4 sold, offered for sale, or distributed in California shall offer users
5 the option to apply difficult to remove provenance data to
6 nonsynthetic content produced by that device.

7 (2) A user shall have the option to not apply provenance data
8 and any other information attached to nonsynthetic content
9 produced by their device and to customize the types of provenance
10 data attached to nonsynthetic content produced by their device,
11 including by removing any personally identifiable information.
12 Personally identifiable information, including geolocation, shall
13 not be included in provenance data by default.

14 (3) Recording devices subject to the requirements of this
15 subdivision shall clearly inform users of the existence of the
16 settings relating to provenance data upon a user's first use of the
17 recording function on the recording device.

18 (4) When a recording device's recording function is in use, the
19 recording device shall contain a clear indicator when provenance
20 data is being applied.

21 (5) The option to apply provenance data to nonsynthetic content
22 produced by a recording device, as described by paragraph (1),
23 shall also be applied to nonsynthetic content produced using
24 third-party applications that bypass default recording applications
25 in order to offer recording functionalities.

26 (6) The application of provenance data shall be compatible with
27 state-of-the-art widely adopted and relevant industry standards.

28 (b) If technically feasible and secure, a recording device
29 manufacturer shall offer a software or firmware update enabling
30 a user of a recording device manufactured before July 1, 2026,
31 and purchased in California to do both of the following:

32 (1) Apply difficult to remove provenance data to the
33 nonsynthetic content created by the device.

34 (2) Decode any provenance data attached to the nonsynthetic
35 content created by the device.

36 22949.90.3. (a) A large online platform capable of
37 disseminating content that would falsely appear to a reasonable
38 person to depict real-life persons, objects, places, entities, or events
39 shall use labels to disclose any machine-readable provenance data
40 detected in synthetic content distributed on its platform.

1 (1) To the extent technically feasible, the labels shall indicate
2 whether provenance data is available.

3 (2) A user shall be able to click or tap on a label to inspect
4 provenance data in an easy-to-understand format.

5 (b) The disclosure required under subdivision (a) shall be readily
6 legible to an average viewer or, if the content is in audio format,
7 shall be clearly audible.

8 (c) If content uploaded to or distributed on a large online
9 platform by a user does not contain provenance data or if the
10 content's provenance data cannot be interpreted or detected by the
11 platform using technically feasible methods, a large online platform
12 shall label the content as having unknown provenance.

13 (d) A large online platform shall add the following provenance
14 data to digital content published on their platform:

15 (1) The name of the platform on which the content was
16 published.

17 (2) The date and time of publishment on the platform.

18 (3) The term "unknown creation process" if the digital content
19 did not contain any previously applied provenance data at the time
20 it was published on the platform.

21 (e) (1) Notwithstanding anything to the contrary in this section,
22 for purposes of labeling and disclosing provenance data of sound
23 recordings and music videos, a large online platform shall use a
24 visual, not an audio, disclosure for sound recordings and music
25 videos that contains all of the following:

26 (A) The artist.

27 (B) The track.

28 (C) The copyright holder or licensor information.

29 (2) A large online platform shall comply with the visual
30 disclosure requirement described in paragraph (1) to the extent
31 that those sound recordings and music videos have not been solely
32 generated by a GenAI system, extended or modified by a GenAI
33 system without the authorization of the copyright holder whose
34 work has been modified or extended, or modified by a GenAI
35 system to imitate or be readily identifiable as another person and
36 that other person has not authorized the modification.

37 (f) This section shall not apply to any product, service, website,
38 or application that provides predominantly non-user-generated
39 video game, television, streaming, or movie experiences.

1 22949.90.4. (a) Beginning July 1, 2026, and annually
2 thereafter, a large online platform shall produce a transparency
3 report that identifies moderation of deceptive synthetic content on
4 their platform.

5 (b) The report required by subdivision (a) may include
6 assessments of the distribution of illegal generative AI-generated
7 child sexual abuse materials, nonconsensual intimate imagery,
8 disinformation related to elections or public health, or other
9 instances where synthetic or potentially deceptive content was
10 identified and removed by the platform.

11 22949.90.5. The Department of Technology may assess an
12 administrative penalty pursuant to the following:

13 (a) If a violation of this chapter is intentional or is the result of
14 grossly negligent conduct, a penalty of up to ~~five~~ *one* hundred
15 thousand dollars ~~(\$500,000)~~ *(\$100,000)* for each violation.

16 (b) If a violation of this chapter is unintentional or is not the
17 result of grossly negligent conduct, a penalty of up to ~~fifty~~
18 *twenty-five* thousand dollars ~~(\$50,000)~~ *(\$25,000)* for each violation.

19 22949.90.6. (a) The Digital Content Provenance Administrative
20 Fund is hereby created in the State Treasury.

21 (b) All penalties collected by the Department of Technology
22 under Section 22949.90.5 shall be deposited in the Digital Content
23 Provenance Administrative Fund.

24 (c) Upon appropriation by the Legislature for this express
25 purpose, moneys in the Digital Content Provenance Administrative
26 Fund may be expended by the Department of Technology to
27 administer this chapter.

28 22949.90.7. This chapter shall become operative on July 1,
29 2026.

30 22949.91. The provisions of this chapter are severable. If any
31 provision of this chapter or its application is held invalid, that
32 invalidity shall not affect other provisions or applications that can
33 be given effect without the invalid provision or application.