

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 Provenance, Authenticity and Watermarking 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.

Beginning February 1, 2025, this bill, the California Provenance, Authenticity and Watermarking Standards Act, would require a generative artificial intelligence (AI) system provider, as defined, to take certain actions to assist in the disclosure of provenance data to mitigate harms caused by inauthentic content, including placing imperceptible and maximally indelible watermarks containing provenance data into content created by an AI system that the generative AI system provider makes available: *provider to, among other things,*

place imperceptible and maximally indelible watermarks containing provenance data into synthetic content produced or significantly modified by a generative AI system that the provider makes available, as those terms are defined. The bill would require, within 24 hours of discovering a vulnerability or failure in ~~an~~ a generative AI system, a generative AI ~~system~~ provider to report the vulnerability or failure to the Department of Technology and to notify other generative AI ~~system~~ providers, as specified. The bill would also require a conversational AI system, as defined, to clearly and prominently disclose to users that the conversational AI system receives synthetic content.

Beginning March 1, 2025, this bill would require a large online platform, as defined, to, among other things, use labels to prominently disclose the provenance data found in watermarks or digital signatures in content distributed to users on its platforms, as specified. The bill would require a large online platform to ~~use state-of-the-art techniques, including, but not limited to, analysis of user behavioral signals indicating usage of synthetic content, to detect and label inauthentic text content that is uploaded or distributed by individual users or networks of users.~~ *require a user that uploads or distributes content on its platform to disclose whether the content is synthetic content, as specified.*

Beginning January 1, 2026, this bill would require newly manufactured digital cameras and recording devices sold, offered for sale, or distributed in California to offer users the option to place an authenticity watermark and provenance watermark in the content produced by that device. The bill would require the authenticity watermark and provenance watermark to be compatible with widely used industry standards, as specified. ~~If a digital camera or recording device purchased in California prior to January 1, 2026, is capable of receiving a software or firmware update that would enable a user to place an authenticity watermark and provenance watermark on the content created by the device, the bill would require the device's manufacturer to offer that update.~~ *standards. If technically feasible, the bill would require a camera and recording device manufacturer, as defined, to offer to a user of a digital camera or recording device purchased in California prior to January 1, 2026, a software or firmware update enabling the user to place an authenticity watermark and provenance watermark on the content created by the device.*

Beginning January 1, 2026, and annually thereafter, this bill would also require ~~specified entities~~ *generative AI providers and large online*

platforms to produce a Risk Assessment and Mitigation Report that assesses the risks posed by synthetic content and the harms that have been or could be caused by synthetic content, and harms caused by synthetic content generated in their systems or hosted on their platforms, as prescribed. The bill would require the report to be audited by qualified, independent auditors who are required to assess and either validate or invalidate the claims made in the report, as specified.

This bill would provide that a violation of its provisions may result in an administrative penalty, assessed by the department, of up to \$1,000,000 or 5% of the violator’s annual global revenue, whichever is greater. The bill would require the department to adopt regulations as necessary to implement and carry out the purposes of this act and to review and update those regulations as needed.

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 produce inauthentic images, audio, video, and*
- 5 *text content 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 GenAI content to be clearly disclosed and*
- 8 *labeled.*
- 9 (c) *Failing to appropriately label GenAI content can skew*
- 10 *election results, enable academic dishonesty, and erode trust in*
- 11 *the online information ecosystem.*
- 12 (d) *The Legislature should act to adopt standards pertaining*
- 13 *to the clear disclosure and labeling of GenAI content, in order to*
- 14 *alleviate harms caused by the misuse of these technologies.*
- 15 (e) *The Legislature should push for the creation of tools that*
- 16 *allow Californians to assess the authenticity of online content.*
- 17 (f) *The Legislature should require online platforms to label*
- 18 *inauthentic content produced by GenAI.*
- 19 (g) *Through these actions, the Legislature can help to ensure*
- 20 *that Californians remain safe and informed.*
- 21 SEC. 2. *Chapter 41 (commencing with Section 22949.90) is*
- 22 *added to Division 8 of the Business and Professions Code, to read:*

1
2 *CHAPTER 41. CALIFORNIA PROVENANCE, AUTHENTICITY, AND*
3 *WATERMARKING STANDARDS*

4
5 22949.90. *For purposes of this chapter, the following*
6 *definitions apply:*

7 (a) *“AI red-teaming” means a structured testing effort to find*
8 *flaws and vulnerabilities in an AI system, including, but not limited*
9 *to, harmful or discriminatory outputs, unforeseen or undesirable*
10 *system behaviors, limitations, or potential risks associated with*
11 *misuse of the system.*

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

17 (c) *“Authentic content” means images, videos, audio, or text*
18 *created by human beings without any modifications or with only*
19 *minor modifications that do not lead to significant changes to the*
20 *perceived contents or meaning of the content. Minor modifications*
21 *include, but are not limited to, changes to brightness or contrast*
22 *of images, removal of background noise in audio, and spelling or*
23 *grammar corrections in text.*

24 (d) *“Conversational AI system” means chatbots and other*
25 *audio- or video-based systems that can hold humanlike*
26 *conversations through digital media, including, but not limited to,*
27 *online calling, phone calling, video conferencing, messaging,*
28 *application or web-based chat interfaces, or other conversational*
29 *interfaces. Conversational AI systems include, but are not limited*
30 *to, chatbots for customer service or entertainment purposes*
31 *embedded in internet websites and applications.*

32 (e) *“Digital signature” means a digital method that allows a*
33 *user to sign a piece of authentic or synthetic content with their*
34 *name or device information, verifying that they created the content.*

35 (f) *“Generative AI hosting platform” means an online repository*
36 *or other internet website that makes generative AI systems*
37 *available for download.*

38 (g) *“Generative AI provider” means an organization or*
39 *individual that creates, codes, substantially modifies, or otherwise*
40 *produces a generative AI system.*

1 (h) “Generative AI system” means an artificial intelligence
2 system that generates derived synthetic content, including images,
3 videos, audio, text, and other digital content.

4 (i) “Inauthentic content” means synthetic content that is so
5 similar to authentic content that it could be mistaken as authentic.

6 (j) “Large online platform” means a public-facing internet
7 website, web application, or digital application, including a social
8 network, video sharing platform, messaging platform, advertising
9 network, or search engine that had at least 1,000,000 California
10 users during the preceding 12 months.

11 (k) “Maximally indelible watermark” means a watermark that
12 is designed to be as difficult to remove as possible using
13 state-of-the-art techniques and relevant industry standards.

14 (l) “Provenance data” means data that identifies the origins of
15 synthetic content, including, but not limited to, the following:

16 (1) The name of the generative AI provider.

17 (2) The name and version number of the AI system that
18 generated the content.

19 (3) The time and date of the creation.

20 (4) The portions of content that are synthetic.

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

24 (n) “Watermark” means information that is embedded into a
25 generative AI system’s output for the purpose of conveying its
26 synthetic nature, identity, provenance, history of modifications,
27 or history of conveyance.

28 (o) “Watermark decoders” means freely available software
29 tools or online services that can read or interpret watermarks and
30 output the provenance data embedded in them.

31 22949.90.1. (a) A generative AI provider shall do all of the
32 following:

33 (1) Place imperceptible and maximally indelible watermarks
34 containing provenance data into synthetic content produced or
35 significantly modified by a generative AI system that the provider
36 makes available.

37 (A) If a sample of synthetic content is too small to contain the
38 required provenance data, the provider shall, at minimum, attempt
39 to embed watermarking information that identifies the content as
40 synthetic and provide the following provenance information in

1 order of priority, with clause (i) being the most important, and
2 clause (iv) being the least important:

3 (i) The name of the generative AI provider.

4 (ii) The name and version number of the AI system that
5 generated the content.

6 (iii) The time and date of the creation of the content.

7 (iv) If applicable, the specific portions of the content that are
8 synthetic.

9 (B) To the greatest extent possible, watermarks shall be designed
10 to retain information that identifies content as synthetic and gives
11 the name of the provider in the event that a sample of synthetic
12 content is corrupted, downscaled, cropped, or otherwise damaged.

13 (2) Develop downloadable watermark decoders that allow a
14 user to determine whether a piece of content was created with the
15 provider's system, and make those tools available to the public.

16 (A) The watermark decoders shall be easy to use by individuals
17 seeking to quickly assess the provenance of a single piece of
18 content.

19 (B) The watermark decoders shall adhere, to the greatest extent
20 possible, to relevant national or international standards.

21 (3) Conduct AI red-teaming exercises involving third-party
22 experts to test whether watermarks can be easily removed from
23 synthetic content produced by the provider's generative AI systems,
24 as well as whether the provider's generative AI systems can be
25 used to falsely add watermarks to otherwise authentic content.
26 Red-teaming exercises shall be conducted before the release of
27 any new generative AI system and annually thereafter.

28 (A) If a provider allows their generative AI systems to be
29 downloaded and modified, the provider shall additionally conduct
30 AI red-teaming to assess whether their systems' watermarking
31 functionalities can be disabled.

32 (B) A provider shall make summaries of its AI red-teaming
33 exercises publicly available in a location linked from the home
34 page of the provider's internet website, using a clearly labeled
35 link that has a similar look, feel, and size relative to other links
36 on the same web page. The provider shall remove from the
37 summaries any details that pose an immediate risk to public safety.

38 (C) A provider shall submit full reports of its AI red-teaming
39 exercises to the Department of Technology within six months of
40 conducting a red-teaming exercise pursuant to this section.

1 **(b)** *A generative AI provider may continue to make available a*
2 *generative AI system that was made available before the date upon*
3 *which this act takes effect and that does not have watermarking*
4 *capabilities as described by paragraph (1) of subdivision (a), if*
5 *either of the following conditions are met:*

6 **(1)** *The provider is able to retroactively create and make*
7 *publicly available a decoder that accurately determines whether*
8 *a given piece of content was produced by the provider’s system*
9 *with at least 99 percent accuracy as measured by an independent*
10 *auditor.*

11 **(2)** *The provider conducts and publishes research to definitively*
12 *demonstrate that the system is not capable of producing inauthentic*
13 *content.*

14 **(c)** *Providers and distributors of software and online services*
15 *shall not make available a system, application, tool, or service*
16 *that is designed to remove watermarks from synthetic content.*

17 **(d)** *Generative AI hosting platforms shall not make available a*
18 *generative AI system that does not place maximally indelible*
19 *watermarks containing provenance data into content created by*
20 *the system.*

21 **(e)** **(1)** *Within 24 hours of discovering a vulnerability or failure*
22 *in a generative AI system, a generative AI provider shall report*
23 *the vulnerability or failure to the Department of Technology.*

24 **(A)** *A provider shall notify other generative AI providers that*
25 *may be affected by similar vulnerabilities or failures in a manner*
26 *that allows the other generative AI provider to harden their own*
27 *AI systems against similar risks.*

28 **(B)** *A provider shall notify affected parties, including, but not*
29 *limited to, online platforms, researchers or users who received*
30 *incorrect results from a watermark decoder, or users who produced*
31 *AI content that contained incorrect or insufficient watermarking*
32 *data. A provider shall not be required to notify an affected party*
33 *whose contact information the provider has not previously collected*
34 *or retained.*

35 **(2)** *A provider shall make any report to the Department of*
36 *Technology under this subdivision publicly available in a location*
37 *linked from the home page of the provider’s internet website with*
38 *a clearly labeled link that has a similar look, feel, and size relative*
39 *to other links on the same web page. If public disclosure of the*

1 report under this subdivision could pose public safety risks, a
2 provider may instead do either of the following:

3 (A) Post a summary disclosure of the reported vulnerability or
4 failure.

5 (B) Delay, for no longer than 30 days, the public disclosure of
6 the report until the public safety risks have been mitigated. If a
7 provider delays public disclosure, they shall document their efforts
8 to resolve the vulnerability or failure as quickly as possible in
9 order to meet the reporting requirements under this subdivision.

10 (f) (1) A conversational AI system shall clearly and prominently
11 disclose to users that the conversational AI system generates
12 synthetic content.

13 (A) In visual interfaces, including, but not limited to, text chats
14 or video calling, a conversational AI system shall place the
15 disclosure required under this subdivision in the interface itself
16 and maintain the disclosure's visibility in a prominent location
17 throughout any interaction with the interface.

18 (B) In audio-only interfaces, including, but not limited to, phone
19 or other voice calling systems, a conversational AI system shall
20 verbally make the disclosure required under this subdivision at
21 the beginning and end of a call.

22 (2) In all conversational interfaces of a conversational AI
23 system, the conversational AI system shall, at the beginning of a
24 user's interaction with the system, obtain a user's affirmative
25 consent acknowledging that the user has been informed that they
26 are interacting with a conversational AI system. A conversational
27 AI system shall obtain a user's affirmative consent prior to
28 beginning the conversation.

29 (3) Disclosures and affirmative consent opportunities shall be
30 made available to a user in the language in which the
31 conversational AI system is communicating with the user.

32 (4) The requirements under this subdivision shall not apply to
33 conversational AI systems that do not produce inauthentic content.

34 (g) This section shall become operative on February 1, 2025.

35 22949.90.2. (a) For purposes of this section, the following
36 definitions apply:

37 (1) "Authenticity watermark" means a watermark of authentic
38 content that includes the name of the device manufacturer.

39 (2) "Camera and recording device manufacturer" means the
40 makers of a device that can record photographic, audio, or video

1 content, including, but not limited to, video and still photography
2 cameras, mobile phones with built-in cameras or microphones,
3 and voice recorders.

4 (3) “Provenance watermark” means a watermark of authentic
5 content that includes details about the content, including, but not
6 limited to, the time and date of production, the name of the user,
7 details about the device, and a digital signature.

8 (b) (1) Beginning January 1, 2026, newly manufactured digital
9 cameras and recording devices sold, offered for sale, or distributed
10 in California shall offer users the option to place an authenticity
11 watermark and provenance watermark in the content produced
12 by that device.

13 (2) A user shall have the option to remove the authenticity and
14 provenance watermarks from the content produced by their device.

15 (3) Authenticity watermarks shall be turned on by default, while
16 provenance watermarks shall be turned off by default.

17 (4) Newly manufactured digital cameras and recording devices
18 subject to the requirements of this subdivision shall clearly inform
19 a user of the existence of the authenticity and provenance
20 watermarks settings upon the user’s first use of the camera or the
21 recording function on the recording device.

22 (A) When a camera or audio recording application is open, a
23 newly manufactured digital camera or recording device shall have
24 a clear indicator that a watermark is being applied.

25 (B) A newly manufactured digital camera or recording device
26 shall allow the user to adjust the watermarks settings.

27 (5) Authenticity and provenance watermarks shall, if enabled,
28 be applied to authentic content produced using third-party
29 applications that bypass default camera or recording applications
30 in order to offer camera or audio recording functionalities.

31 (c) The authenticity watermark and provenance watermark, as
32 required by subdivision (b), shall be compatible with widely used
33 industry standards.

34 (d) Beginning January 1, 2026, a camera and recording device
35 manufacturer shall offer a software or firmware update enabling
36 a user to place an authenticity watermark and provenance
37 watermark on the content created by the device to a user of a
38 digital camera or recording device purchased in California prior
39 to January 1, 2026, if technically feasible.

1 22949.90.3. (a) Beginning March 1, 2025, a large online
2 platform shall use labels to prominently disclose the provenance
3 data found in watermarks or digital signatures in content
4 distributed to users on its platforms.

5 (1) The labels shall indicate whether content is fully synthetic,
6 partially synthetic, authentic, authentic with minor modifications,
7 or does not contain a watermark.

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

10 (b) The disclosure required under subdivision (a) shall be
11 readily legible to an average viewer or, if the content is in audio
12 format, shall be clearly audible. A disclosure in audio content
13 shall occur at the beginning and end of a piece of content and shall
14 be presented in a prominent manner and at a comparable volume
15 and speaking cadence as other spoken words in the content. A
16 disclosure in video content should be legible for the full duration
17 of the video.

18 (c) A large online platform shall use state-of-the-art techniques
19 to detect and label synthetic content that has had watermarks
20 removed or that was produced by generative AI systems without
21 watermarking functionality.

22 (d) (1) A large online platform shall require a user that uploads
23 or distributes content on its platform to disclose whether the
24 content is synthetic content.

25 (2) A large online platform shall include prominent warnings
26 to users that uploading or distributing synthetic content without
27 disclosing that it is synthetic content may result in disciplinary
28 action.

29 (3) A large online platform may provide users with an option
30 to indicate that the user is uncertain whether the content they are
31 uploading or distributing is synthetic content. If a user uploads or
32 distributes content and indicates that they are uncertain of whether
33 the content is synthetic content, a large online platform shall
34 indicate that the uploaded or distributed content is possibly
35 synthetic and shall display that indication in a manner that is
36 visible or audible to viewers or listeners of the content.

37 (e) A large online platform shall use state-of-the-art techniques
38 to detect and label text-based inauthentic content that is uploaded
39 by users.

1 (f) A large online platform shall make accessible a verification
2 process for users to apply a digital signature to authentic content.
3 The verification process shall include options that do not require
4 disclosure of personal identifiable information.

5 22949.90.4. (a) (1) Beginning January 1, 2026, and annually
6 thereafter, generative AI providers and large online platforms
7 shall produce a Risk Assessment and Mitigation Report that
8 assesses the risks posed and harms caused by synthetic content
9 generated by their systems or hosted on their platforms.

10 (2) The report shall include, but not be limited to, assessments
11 of the distribution of AI-generated child sexual abuse materials,
12 nonconsensual intimate imagery, disinformation related to
13 elections or public health, plagiarism, or other instances where
14 synthetic or inauthentic content caused or may have the potential
15 to cause harm.

16 (b) The report required under subdivision (a) shall be audited
17 by qualified, independent auditors who shall assess and either
18 validate or invalidate the claims made in the report. Auditors shall
19 use state-of-the-art techniques to assess reports, and shall adhere
20 to relevant national and international standards.

21 22949.90.5. A violation of this chapter may result in an
22 administrative penalty, assessed by the Department of Technology,
23 of up to one million dollars (\$1,000,000) or 5 percent of the
24 violator's annual global revenue, whichever is greater.

25 22949.90.6. Within 90 days of the date upon which this act
26 takes effect, the Department of Technology shall adopt regulations
27 to implement and carry out the purposes of this chapter. The
28 department shall review and update its regulations relating to the
29 implementation of this chapter as needed, including, but not limited
30 to, adopting specific national or international standards for
31 provenance, authenticity, watermarking, and digital signatures,
32 as long as the standards do not weaken the provisions of this
33 chapter.

34 22949.91. The provisions of this chapter are severable. If any
35 provision of this chapter or its application is held invalid, that
36 invalidity shall not affect other provisions or applications that can
37 be given effect without the invalid provision or application.

38 ~~SECTION 1. The Legislature finds and declares all of the~~
39 ~~following:~~

1 (a) In light of the widespread adoption of generative artificial
2 intelligence (AI) technologies that are increasingly able to generate
3 inauthentic images, audio, video, and text content, sometimes
4 called “deepfakes,” it is increasingly important that the provenance
5 of this content be clearly disclosed and labeled in ways that could
6 both prevent harms, reduce the severity of harms, and also make
7 it more difficult and costly for bad actors to cause these harms.

8 (b) (1) The harms caused by inauthentic content that is
9 presented as authentic span a wide gamut of fields. These harms
10 can be costly and deeply damaging to individuals and society.
11 Some prominent categories of harm include scams and fraud, child
12 sexual abuse material (CSAM) and nonconsensual intimate imagery
13 (NCII), disinformation, and plagiarism and academic integrity.

14 (A) The Federal Trade Commission reports that people lost \$8.8
15 billion to scams in 2022. This is only expected to increase with
16 broad access to generative AI tools, with one study already finding
17 a 1,265-percent increase in phishing scam emails since the fourth
18 quarter of 2022.

19 (B) A recent study found that a version of one of the most
20 popular AI image generator tools was trained on CSAM, which
21 facilitates production in derivative models of both CSAM and
22 NCII of children. Another recent study found 34 different AI
23 “undressing” tools that produce realistic naked images based on
24 clothed photographs of individuals. Many of these tools only work
25 on women. One recent case involved a high school where 30
26 different teenage girls were victims of these tools.

27 (C) From politics to public safety, disinformation can have a
28 multiplicity of negative impacts on society. It can skew election
29 results, as may have been the case with a convincing audio
30 deepfake in Slovakia. In the lead up to the 2024 United States
31 elections, we have already seen candidate- and party-produced
32 advertisements using deepfake audio and video content. A deepfake
33 of an explosion at the Pentagon caused a brief but significant stock
34 market dip.

35 (D) Educational institutions have struggled to adopt clear best
36 practices for assigning students at-home writing assignments since
37 the widespread use of high-quality AI text generation tools.

38 (2) More broadly, erosion of trust in the information ecosystem,
39 sometimes referred to as “truth decay,” can be shown to increase
40 polarization and stands to be greatly exacerbated by our ongoing

1 failure to establish norms for clearly disclosing and labeling the
2 provenance of digital media:

3 (e) (1) For the above-described reasons, this act provides for
4 the phased introduction of the California Provenance, Authenticity
5 and Watermarking Standards (PAWS), which aim to provide the
6 public with tools to understand how the content that they see across
7 digital media was produced and if it is, in fact, authentic.

8 (2) The PAWS require all producers of media to embed
9 maximally indelible and privacy-preserving content provenance
10 data into the content that they generate, whether AI-generated or
11 authentic. The PAWS also require all large online platforms to
12 display clearly understandable labels on content that alert users to
13 its provenance, or to the absence of available provenance data.

14 (3) The provenance data required by PAWS is narrowly tailored
15 to provide consumers with factual information about the
16 authenticity or inauthenticity of images, audio, video, or text
17 content in order to prevent consumer deception.

18 (4) While the PAWS alone will not fully eliminate the harms
19 described above, it has the capacity to dramatically reduce these
20 harms by signaling to the would-be audiences for and victims of
21 these harms that inauthentic content will not be easily mistaken
22 for authentic content. The PAWS will also significantly reduce
23 the burden on large online platforms that would like to display
24 content provenance data to their users, but are not able to do so
25 because of the lack of an industrywide standard for embedding of
26 provenance data in their content.

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

29

30 CHAPTER 41. CALIFORNIA PROVENANCE, AUTHENTICITY AND
31 WATERMARKING STANDARDS

32

33 22949.90. For purposes of this chapter, the following
34 definitions apply:

35 (a) “AI red-teaming” means a structured testing effort to find
36 flaws and vulnerabilities in an AI system, often in a controlled
37 environment and in collaboration with developers of artificial
38 intelligence (AI) and is most often performed by dedicated “red
39 teams” that adopt adversarial methods to identify flaws and
40 vulnerabilities, including, but not limited to, harmful or

1 discriminatory outputs from an AI system, unforeseen or
2 undesirable system behaviors, limitations, or potential risks
3 associated with misuse of the system.

4 (b) “AI system” means a machine-based system that, for explicit
5 or implicit objectives, infers, from the input it receives, how to
6 generate outputs, including, but not limited to, predictions, content,
7 recommendations, or decisions that may influence physical or
8 virtual environments.

9 (c) “Authentic content” means images, videos, audio clips, or
10 text created by human beings without any modifications or with
11 only minor modifications that do not lead to significant changes
12 to the perceived contents or meaning of the content. Minor
13 modifications include, but are not limited to, changes to brightness
14 or contrast of images, removal of background noise in audio, and
15 spelling or grammar corrections in text.

16 (d) “Camera and recording device manufacturers” means the
17 makers of a device that can record photographs, audio, or video
18 content, including, but not limited to, video and still photography
19 cameras, mobile phones with built-in cameras or microphones,
20 and voice recorders.

21 (e) “Conversational AI systems” means chatbots and other
22 audio- or video-based systems that can hold humanlike
23 conversations through digital media, including, but not limited to,
24 online calling, phone calling, video conferencing, messaging,
25 application or web-based chat interfaces, or other conversational
26 interfaces. Conversational AI systems includes chatbots for
27 customer service or entertainment purposes embedded in internet
28 websites and applications.

29 (f) “Digital signature” means a digital method that allows a user
30 to sign a piece of authentic or synthetic content with their name
31 or device information, verifying that they created the content, and
32 that is included in the content’s provenance data.

33 (g) “Generative AI system” means the class of AI models that
34 emulates the structure and characteristics of input data to generate
35 derived synthetic content, including images, videos, audio, text,
36 and other digital content. The synthetic content generated may,
37 but does not necessarily have to, be of the same content type as
38 the input data.

39 (h) “Generative AI system distributors” means organizations
40 or individuals that distribute generative AI systems, or substantial

1 components thereof, such as model weights, in ways that can be
2 downloaded and used by individuals locally on their own hardware,
3 or modified or incorporated into other products or services.

4 (i) “Generative AI system providers” means organizations or
5 individuals that make AI systems, or substantial components
6 thereof, available on the market, put them into service, provide
7 them as standalone models or embed them in other systems or
8 products, or provide them under free and open source licenses as
9 a service, or through other large online platforms. Generative AI
10 system distributors include repositories or hosting internet websites
11 that make AI systems available for download, even if those
12 repositories are not the original makers of the AI systems they
13 make available, and providers of conversational AI systems.

14 (j) “Inauthentic content” means synthetic content that is so
15 similar to authentic content that it could be mistaken as authentic.

16 (k) “Large online platform” means a public-facing internet
17 website, web application, or digital application, including a social
18 network, video sharing platform, messaging platform, advertising
19 network, or search engine that had at least 1,000,000 California
20 users during the preceding 12 months.

21 (l) “Legacy generative AI systems” means generative AI systems
22 created and released before the date upon which this act takes
23 effect.

24 (m) “Maximally indelible watermarks” means watermarks that
25 are as difficult to remove as is possible with currently available
26 techniques and that are tested through AI red-teaming.

27 (n) “Provenance data” means data that includes the name of the
28 AI system that generated the content, the underlying AI models
29 that were part of the AI system, the time and date of the creation
30 of the content, and, if applicable, which specific portions of the
31 content are synthetic.

32 (o) “Synthetic content” means information, including images,
33 videos, audio clips, and text, that has been significantly modified
34 or generated by algorithms, including by AI.

35 (p) “Tamper-evident” means a type of watermark that contains
36 provenance data that cannot be modified without leaving evidence
37 of tampering.

38 (q) “Watermark” means embedded information, typically
39 difficult to remove, into outputs created, including into
40 photographs, videos, audio clips, or text, for the purposes of

1 verifying the authenticity of the output or the identity or
2 characteristics of its provenance, modifications, or conveyance.

3 (r) “Watermark decoders” means freely available software tools
4 or online services that can read or interpret watermarks and output
5 the provenance data embedded in them.

6 ~~22949.90.1. (a) A generative AI system provider shall do all~~
7 ~~of the following:~~

8 ~~(1) Place imperceptible and maximally indelible watermarks~~
9 ~~containing provenance data into content created by an AI system~~
10 ~~that the generative AI system provider makes available. If the~~
11 ~~content is manipulated in ways that are intended to remove a~~
12 ~~watermark, or if a content sample, including a cropped portion of~~
13 ~~an image, a truncated segment of a generated piece of audio or~~
14 ~~video content, or a small amount of text content, is damaged or~~
15 ~~becomes too small to contain the required provenance data, the~~
16 ~~generative AI system provider shall, at minimum, embed~~
17 ~~watermarking information that identifies the content as synthetic~~
18 ~~and gives the name of the generative AI system provider.~~

19 ~~(2) Provide downloadable software tools or online services to~~
20 ~~determine whether a piece of content was created with the~~
21 ~~provider’s system and make those tools available to all large online~~
22 ~~platforms and the public.~~

23 ~~(A) A generative AI system provider shall produce the software~~
24 ~~tools or online services in ways that are easy to use manually by~~
25 ~~individuals seeking to quickly assess the provenance of a single~~
26 ~~piece of content and to use in an efficient and automated fashion~~
27 ~~by online platforms or researchers seeking to assess the provenance~~
28 ~~of hundreds or thousands of pieces of content per minute.~~

29 ~~(B) A generative AI system provider shall produce the software~~
30 ~~tools or online services to be interoperable to the greatest extent~~
31 ~~possible with decoders made available by other providers or~~
32 ~~organizations. The decoders shall adhere, to the greatest extent~~
33 ~~possible, to relevant national or international standards, if available.~~

34 ~~(3) Conduct AI red-teaming exercises that involve third-party~~
35 ~~experts to test whether watermarks can be easily removed from~~
36 ~~their synthetic content and whether AI systems could be used to~~
37 ~~falsely add watermarks to authentic content that indicates that it~~
38 ~~is inauthentic.~~

39 ~~(A) If a generative AI system provider intends to allow their~~
40 ~~systems to be downloaded and modified, the provider shall conduct~~

1 ~~AI red-teaming to assess whether the systems' watermarking~~
2 ~~functionalities can be disabled to generate deceptive, inauthentic~~
3 ~~content.~~

4 ~~(B) A generative AI system provider shall make summaries of~~
5 ~~its AI red-teaming exercises publicly available in a location linked~~
6 ~~from the home page of the generative AI system provider's internet~~
7 ~~website with a clearly labeled link that has a similar look, feel, and~~
8 ~~size relative to other links on the same web page and shall remove~~
9 ~~from the summaries details that may pose immediate public~~
10 ~~security. A generative AI system provider shall submit full reports~~
11 ~~of its AI red-teaming exercises to the Department of Technology~~
12 ~~within six months from the date upon which this act takes effect,~~
13 ~~and annually thereafter.~~

14 ~~(b) A generative AI system provider may continue providing~~
15 ~~legacy generative AI systems that do not have watermarking~~
16 ~~capabilities as described by paragraph (1) of subdivision (a) if~~
17 ~~either of the following applies:~~

18 ~~(1) The generative AI system provider is able to retrospectively~~
19 ~~create and make publicly available a decoder that accurately~~
20 ~~determines whether content generated by the legacy generative AI~~
21 ~~system is synthetic with at least 99 percent accuracy as measured~~
22 ~~by an independent auditor.~~

23 ~~(2) The generative AI system provider conducts and publishes~~
24 ~~research to definitively demonstrate that the legacy generative AI~~
25 ~~system is not capable of producing inauthentic content of a quality~~
26 ~~that is sufficiently realistic to be potentially mistaken for authentic~~
27 ~~content.~~

28 ~~(c) A generative AI system provider, generative AI system~~
29 ~~distributor, or other provider, or a distributor of software or online~~
30 ~~services shall not make available or distribute an AI system,~~
31 ~~application, tool, or service that has the capacity to remove~~
32 ~~watermarks from synthetic content or an AI system that has the~~
33 ~~capacity to remove watermarks from synthetic content but was not~~
34 ~~explicitly designed for that purpose.~~

35 ~~(d) (1) If a generative AI system provider distributes a~~
36 ~~generative AI system in a way that allows the generative AI system~~
37 ~~to be modified by others, the generative AI system provider shall~~
38 ~~ensure that the generative AI system does not allow for removal~~
39 ~~of the system's watermarking functionality.~~

1 ~~(2) A generative AI system provider shall not distribute a~~
2 ~~generative AI system or make a generative AI system available~~
3 ~~for use if the system's watermarking functionality can be removed~~
4 ~~by others.~~

5 ~~(e) (1) Within 24 hours of discovering a vulnerability or failure~~
6 ~~in an AI system, a generative AI system provider shall report the~~
7 ~~vulnerability or failure to the Department of Technology. A~~
8 ~~generative AI system provider shall also notify other generative~~
9 ~~AI system providers that may be affected by similar vulnerabilities~~
10 ~~or failures in a manner that allows the other generative AI system~~
11 ~~provider to harden their own AI systems against similar risks.~~

12 ~~(2) A generative AI system provider shall also notify any~~
13 ~~affected party. Affected parties include, but are not limited to, an~~
14 ~~online platform, researchers or users who received incorrect results~~
15 ~~from a watermark decoder, or users who produced AI content that~~
16 ~~contained incorrect or insufficient watermarking data.~~

17 ~~(3) A generative AI system provider shall make any report to~~
18 ~~the Department of Technology under this subdivision publicly~~
19 ~~available in a location linked from the home page of the generative~~
20 ~~AI system provider's internet website with a clearly labeled link~~
21 ~~that has a similar look, feel, and size relative to other links on the~~
22 ~~same web page.~~

23 ~~(4) If public disclosure of the report under this subdivision could~~
24 ~~pose public safety risks, a generative AI system provider may~~
25 ~~instead do either of the following:~~

26 ~~(A) Post a summary disclosure of the reported vulnerability or~~
27 ~~failure.~~

28 ~~(B) Delay, for no longer than 30 days, the public disclosure of~~
29 ~~the report until the public safety risks have been mitigated. If a~~
30 ~~generative AI system provider delays public disclosure, the~~
31 ~~provider shall document and demonstrate that they moved as~~
32 ~~quickly as possible to resolve the vulnerability or failure and meet~~
33 ~~the reporting requirements under this subdivision.~~

34 ~~(f) (1) A conversational AI system shall clearly and prominently~~
35 ~~disclose to users that the conversational AI system receives~~
36 ~~synthetic content.~~

37 ~~(A) In visual interfaces, including, but not limited to, text chats~~
38 ~~or video calling, a conversational AI system shall place the~~
39 ~~disclosure required under this subdivision in the interface itself~~

1 and maintain the disclosure’s visibility in a prominent location
2 throughout any interaction with the interface.

3 (B) In audio-only interfaces, including, but not limited to, phone
4 or other voice calling systems, a conversational AI system shall
5 verbally make the disclosure required under this subdivision at the
6 beginning and end of a call.

7 (2) In all conversational interfaces of a conversational AI system,
8 the conversational AI system shall, at the beginning of a user’s
9 interaction with the system, obtain a user’s affirmative consent
10 acknowledging that the user has been informed that they are
11 interacting with a conversational AI system. A conversational AI
12 system shall obtain a user’s affirmative consent prior to beginning
13 the conversation.

14 (3) Disclosures and affirmative consent opportunities shall be
15 made available to a user in the language in which the
16 conversational AI system is communicating with the user.

17 (4) The requirements under this subdivision shall not apply to
18 conversational AI systems that produce content that could not be
19 reasonably mistaken as authentic.

20 (g) This section shall become operative on February 1, 2025.

21 22949.90.2. (a) For purposes of this section, the following
22 definitions apply:

23 (1) “Authenticity watermark” means a watermark of authentic
24 content that includes the name of the device manufacturer.

25 (2) “Provenance watermark” means a watermark of authentic
26 content that includes details about the content, including, but not
27 limited to, the time and date of production, the name of the user,
28 details about the device, and a digital signature.

29 (b) (1) Beginning January 1, 2026, newly manufactured digital
30 cameras and recording devices sold, offered for sale, or distributed
31 in California shall offer users the option to place an authenticity
32 watermark and provenance watermark in the content produced by
33 that device.

34 (2) A user shall have the option to remove the authenticity and
35 provenance watermarks from the content produced by their device.

36 (3) Authenticity watermarks shall be turned on by default, while
37 provenance watermarks shall be turned off by default.

38 (4) Newly manufactured digital cameras and recording devices
39 subject to the requirements of this subdivision shall clearly inform
40 a user of the existence of the authenticity and provenance

1 watermarks settings upon the user's first use of the camera or the
2 recording function on the recording device.

3 (A) When a camera or audio recording application is open, a
4 newly manufactured digital camera or recording device shall have
5 a clear indicator that a watermark is being applied.

6 (B) A newly manufactured digital camera or recording device
7 shall allow the user to adjust the watermarks settings.

8 (5) The authenticity watermark and provenance watermark shall
9 persist in content produced using newly manufactured digital
10 cameras or recording devices even if the content is created or
11 recorded in third-party applications that offer camera or audio
12 recording functionalities.

13 (e) The authenticity watermark and provenance watermark, as
14 required in subdivision (b), shall be compatible with widely used
15 industry standards, including the Coalition for Content Provenance
16 and Authenticity's content credentials.

17 (d) Beginning January 1, 2026, a camera and recording device
18 manufacturer shall offer a software or firmware update enabling
19 a user to place an authenticity watermark and provenance
20 watermark on the content created by the device to a user of a digital
21 camera or recording device purchased in California prior to January
22 1, 2026. This requirement shall only apply if a digital camera or
23 recording device purchased in California prior to January 1, 2026,
24 is capable of receiving a software or firmware update that would
25 enable the user to place an authenticity watermark and provenance
26 watermark on the content created by the device.

27 22949.90.3. (a) Beginning March 1, 2025, a large online
28 platform shall use labels to prominently disclose the provenance
29 data found in watermarks or digital signatures in content distributed
30 to users on its platforms by making use of the data contained in
31 watermarks and digital signatures embedded in content using
32 widely used industry standards, including the Coalition for Content
33 Provenance and Authenticity's content credentials, and synthetic
34 content decoders provided by generative AI system providers.

35 (1) The labels shall indicate whether content is fully synthetic,
36 partially synthetic, authentic, authentic with minor modifications,
37 or does not contain a watermark.

38 (2) A user shall be allowed to click or tap on a label to inspect
39 all available provenance data in an easy-to-understand format.

1 ~~(b) The disclosure required under subdivision (a) shall be readily~~
2 ~~legible to an average viewer or, if the content is in audio format,~~
3 ~~shall be clearly audible. A disclosure in audio and video content~~
4 ~~shall occur at the beginning and end of a piece of content and shall~~
5 ~~be presented in a prominent manner and at a comparable volume~~
6 ~~and speaking cadence as other spoken words in the content.~~

7 ~~(c) A large online platform shall use state-of-the-art techniques~~
8 ~~to detect and label synthetic content that has had watermarks~~
9 ~~removed or was produced by AI systems without watermarking~~
10 ~~functionality.~~

11 ~~(d) (1) A large online platform shall require a user that uploads~~
12 ~~or distributes content on its platform to disclose whether the~~
13 ~~uploaded or distributed content is synthetic content.~~

14 ~~(2) A large online platform shall include prominent warnings~~
15 ~~to users that uploading or distributing synthetic content without~~
16 ~~disclosing that it is synthetic content is not permissible and will~~
17 ~~result in disciplinary action by the large online platform.~~

18 ~~(3) A large online platform may provide users with an option~~
19 ~~to indicate that the user is uncertain whether the content they are~~
20 ~~uploading or distributing is synthetic content. If a user uploads or~~
21 ~~distributes content and indicates that they are uncertain of whether~~
22 ~~the content is synthetic content, a large online platform shall~~
23 ~~indicate that the uploaded or distributed content is possibly~~
24 ~~synthetic and shall display that indication in a manner that is visible~~
25 ~~or audible to viewers or listeners of the content.~~

26 ~~(e) (1) A large online platform shall use state-of-the-art~~
27 ~~techniques to detect and label inauthentic text content that is~~
28 ~~uploaded or distributed by individual users or networks of users.~~

29 ~~(2) A large online platform may use a variety of methods to~~
30 ~~detect inauthentic text content, including, but not limited to, the~~
31 ~~following:~~

32 ~~(A) Bulk analysis of collected text content from users or~~
33 ~~networks of users.~~

34 ~~(B) Analysis of user behavioral signals indicating usage of~~
35 ~~synthetic content.~~

36 ~~(C) Assessing large quantities of text generated by users or~~
37 ~~networks of users for watermarked content.~~

38 ~~(D) Considering whether a user's typing cadence indicates~~
39 ~~authenticity or automation.~~

1 ~~(E) Verification that users are matched to unique device~~
2 ~~identifications such as a subscriber identity module (SIM) card,~~
3 ~~international mobile equipment identity (IMEI), or multifactor~~
4 ~~authentication (MFA):~~

5 ~~(3) A large online platform shall also consider account age,~~
6 ~~login frequency, connection to other identity verification services,~~
7 ~~frequency of content uploading or distributing, authenticity of~~
8 ~~original media content, and other on-platform behaviors by a user~~
9 ~~that could be used to detect undisclosed inauthentic content~~
10 ~~production.~~

11 ~~(4) If a large online platform discovers that a user has uploaded~~
12 ~~or distributed inauthentic content and the user did not disclose that~~
13 ~~the uploaded or distributed content is synthetic content pursuant~~
14 ~~to subdivision (d), the large online platform shall disable the~~
15 ~~account of the user that uploaded or distributed the undisclosed~~
16 ~~inauthentic content.~~

17 ~~(f) A large online platform shall make accessible a verification~~
18 ~~process for users to apply a digital signature to content created by~~
19 ~~a human being. The verification process shall include options to~~
20 ~~verify in a variety of methods that do not necessarily require~~
21 ~~disclosure of personal identifiable information, including, but not~~
22 ~~limited to, uploading a government-issued identification and~~
23 ~~matching picture identification or verifying that a user possesses~~
24 ~~a unique device with a SIM card and active phone number.~~

25 ~~22949.90.4. (a) (1) Beginning January 1, 2026, and annually~~
26 ~~thereafter, generative AI system providers, generative AI system~~
27 ~~distributors, and large online platforms shall produce a Risk~~
28 ~~Assessment and Mitigation Report that assesses the risks posed~~
29 ~~by synthetic content and the harms that have been or could be~~
30 ~~caused by synthetic content.~~

31 ~~(2) The report shall include, but is not limited to, assessments~~
32 ~~on the distribution of AI-generated child sexual abuse materials,~~
33 ~~nonconsensual intimate imagery, disinformation related to elections~~
34 ~~or public health, plagiarism, or other instances where synthetic or~~
35 ~~inauthentic content caused or may have the potential to cause harm.~~

36 ~~(3) The report shall be audited by qualified, independent auditors~~
37 ~~who shall assess and either validate or invalidate the claims made~~
38 ~~in the report. An auditor shall assess the report by using~~
39 ~~state-of-the-art techniques and adhering to national and~~

1 international standards for the auditing of AI systems as they
2 become available.

3 (b) The Department of Technology may authorize independent
4 researchers associated with educational institutions or civil society
5 organizations approved by the Department of Technology to access
6 special researcher tools designed to facilitate large-scale and
7 efficient analysis of content and application programming
8 interfaces (APIs) from generative AI system providers and large
9 online platforms for the purposes of generating test content,
10 studying the efficacy of labeling and effects on users, and
11 evaluating the overall effectiveness of this chapter in preventing
12 harms caused by inauthentic content.

13 22949.90.5. A violation of this chapter may result in an
14 administrative penalty, assessed by the Department of Technology,
15 of up to one million dollars (\$1,000,000) or 5 percent of the
16 violator's annual global revenue, whichever is greater.

17 22949.90.6. Within 90 days of the date upon which this act
18 takes effect, the Department of Technology shall adopt regulations
19 as necessary to implement and carry out the purposes of this
20 chapter. The department shall review and update its regulations
21 relating to the implementation of this chapter as needed, including,
22 but not limited to, adopting specific national or international
23 standards for provenance, authenticity, watermarking, and digital
24 signatures, so long as the standards do not weaken the provisions
25 of this chapter.

26 22949.91. The provisions of this chapter are severable. If any
27 provision of this chapter or its application is held invalid, that
28 invalidity shall not affect other provisions or applications that can
29 be given effect without the invalid provision or application.