



**Select Committee on Adopting Artificial Intelligence Hearing
16 August 2024**

Responses to Questions Taken on Notice - Part II

20 September 2024

By Email: aicommitee.sen@aph.gov.au

Part II Questions

1. What is Google's view on what did and didn't work in the South Korean approach to regulating deepfakes around its recent elections?

Our experience in South Korea confirmed for us the importance of maintaining a careful balance between freedom of expression and information, and removal of harmful information, particularly with respect to definitions of what is considered to be 'deepfake' material. It is important that political speech be enabled, and that regulatory measures designed to prevent harm from deepfakes do not also hinder election campaigning. In this context, it is important to focus any policy or regulatory approach on situations where there is clearly an intent to mislead or deceive through the use of AI.

2. What is Google's view on other measures taken by Governments to regulate AI-generated material around elections and in a political context more broadly?

Google has not waited for Governments to regulate the use of AI-generated material in connection with elections. We know our users depend on us to provide reliable and up-to-date information on topics like current candidates, voting processes and election results — and this new technology can make mistakes as it learns or as news breaks.

Out of an abundance of caution, we restrict the types of election-related queries that our own AI tools, such as the Gemini app and web experience, will respond to. As we integrate Gen AI into more consumer experiences, we're also applying election-related restrictions to many of these products, including YouTube AI-generated summaries for Live Chat, Gems, and image generation in Gemini. For many of these queries on Gemini, we also provide a link connecting users directly to Google Search for the latest and most accurate information.

In addition, we were the first tech company to [require](#) election advertisers to prominently disclose when their ads include realistic synthetic content that's been digitally altered or generated, including by AI tools. Our own AI image generation tools are not available in sensitive ad verticals like "political": [please see this Ads Help Centre page](#) for more information.

On Search, our users are able to access our "[About this result](#)" information literacy feature, which gives them context and more information about a source before they click through.

And our users can click on "[About this Image](#)" to see information like when and where similar images first appeared, and where else the image has been seen online, including on news, fact-checking and social media sites. This will provide users with helpful context to determine whether what they are seeing is reliable.

We require creators on YouTube to [disclose](#) when they upload realistic altered or synthetic content and display labels on these videos so viewers know that it contains altered or synthetic content. For election-related content, a label will appear both on the video player and in the video description.

If people ask Gemini questions, they'll see a "[double check](#)" feature that gives them helpful context about the response they are seeing, so they can click through to Google Search and evaluate whether there are authoritative sources supporting what they're seeing.

For content created by our own AI generative tools, we created [SynthID](#), which directly embeds a digital watermarking into text, images, video, and audio so people can identify it as AI-generated.

3. Please advise separately for each of the following Google products and services, whether user data has been used to train Gemini or other Google AI products and services, specifying the extent to which that data has been used, the specific models that have been trained on this data, and how consent was obtained (including whether it was on an opt in, opt out, or compulsory basis):

- a. Google Search**
- b. Google Chrome**
- c. Google Cloud**
- d. Google Drive**

- e. Google Docs Editors suite of products including Docs, Sheets, Slides, Drawings, Forms, Sites, Keep and
- f. YouTube
- g. YouTube Music
- h. Gmail
- i. Pixel smartphone devices
- j. Chromebook devices
- k. Fitbit devices
- l. Pixel Watch devices
- m. Google Maps
- n. Google Earth
- o. Google Street View
- p. Waze
- q. Google Books
- r. Google Voice
- s. Google Voice Search
- t. Google Chat
- u. Google Meet
- v. Google Translate
- w. Google Images
- x. Google Pay & Google Wallet
- y. Google Workspace
- z. Google Assistant
- aa. Google Fit
- bb. Google Photos

Google's AI models are trained primarily on publicly available, crawlable data from the open internet. Such data typically comes from a wide range of sources such as web documents, and code, and often include image, audio, and video data along with text.

Our policies make it clear that publicly available data may be used to help train our AI and language models.

4. Google also serves as a publisher for third party creators through various products. Please advise separately for each of the following Google products and services, whether content uploaded or accessible through the platforms has been used to train Gemini or other Google AI products and services, specifying the extent to which that data has been used, the specific models that have been trained on this data, and how consent was obtained (including whether it was on an opt in, opt out, or compulsory basis):

- a. Google Books
- b. YouTube
- c. Google Scholar
- d. Google Play

Google's AI models are trained primarily on publicly available, crawlable data from the open internet. Such data typically come from a wide range of sources such as web documents, and code, and often include image, audio, and video data along with text.

Our policies make it clear that publicly available data may be used to help train our AI and language models.

5. At the hearing Ms Longcroft said

“What we are doing is legal. We comply with applicable copyright laws in the jurisdictions in which we operate.

Has Google sought legal advice and/or formed an opinion of the legality of the following actions under Australian laws?

Please advise separately for each, and if the answer is yes, please also advise what that advice or opinion is:

a. Google training an AI model in Australia on Australian copyrighted content accessed through Google-owned publishing platforms such as YouTube, Google Books or Google Scholar

b. Google training an AI model overseas on Australian copyrighted content accessed through Google-owned publishing platforms such as YouTube, Google Books or Google Scholar

c. Google or a user of Google using a Google generative AI product such as Gemini, in Australia, to produce an imitation or copy of an Australian copyrighted work, where the model has been trained on that work.

d. Google or a user of Google using a Google generative AI product such as Gemini, while overseas, to produce an imitation or copy of an Australian copyrighted work, where the model has been trained on that work.

e. Google training an AI model in Australia on photos of Australians it has obtained through scraping of “publicly available” web content.

f. Google or a user of Google using a Google generative AI product such as Gemini, while in Australia, to produce a deepfake image of an individual in Australia (intentionally or unintentionally), where the model has been trained on photos of that individual.

To the extent Google has obtained legal advice on these topics that advice is legally privileged and confidential.

For completeness, Google follows the laws of the countries in which AI model training occurs and also offers the Google-Extended opt-out that is honoured irrespective of the applicable local law.

6. At the hearing Ms Longcroft said “our models are trained on publicly available information that is crawled on the open web.” Please advise, answering separately for each:

a. How does google define “publicly available information” in this context?

Publicly available information on the open web is information online that is crawlable in accordance with a website’s robots.txt choices. Such data typically come from a wide range of sources such as web documents, and code, and often include image, audio, and video data along with text.

Before training begins, the data is decomposed in ways that will enable the model to perceive patterns. For example, language-based data is tokenised, meaning sentences are disaggregated into words and portions of words. These tokens are what the model uses to learn about patterns in language, and how to predict the next most-likely token in a sequence.

b. Does publicly available data include content on social media pages that are not hidden by privacy settings? If not, how does Google exclude this data?

If a social media website provides privacy settings that limit the public availability of the content, then it would not be accessible to Google’s web crawlers.

c. Does publicly available data include data that is copyrighted, trademarked or patented? If not, how does Google exclude this data?

Yes.

d. Does publicly available data include pirated content? If not, how does Google exclude this data?

Google takes steps to exclude content that has been reported to it as infringing. Our efforts to fight piracy on the Internet broadly are detailed [online](#) and include our investments in streamlining the copyright removal process for Search results.

7. At the hearing Ms Longcroft was asked whether “social media platforms like Facebook, LinkedIn and Instagram” count as public sources for the purpose for Google web scraping and AI training, to which Ms Longcroft responded: “If it’s publicly available, it’s information; and, of course, within the terms and services of each of those private platforms, it is a question for them to set, with regard to their own users, and to ensure that those users are aware of the purposes for which that information may be used.”

So to clarify – a photo of someone’s children on their Facebook or Instagram account, if that photo has been set to a privacy setting of “Friends only” or “only me”, will be used to train Gemini and other Google generative AI services?

As of this writing, Facebook has opted out via Google-Extended in its robots.txt file, thus no information from crawling Facebook is used to help improve Gemini Apps and Vertex AI generative APIs, including future generations of models that power those products.

8. At the hearing, Google was asked about the need for greater transparency around AI models, to which Ms Doshi responded that “we need to always make sure that we're balancing the needs and privacy of our users and also recognising the importance of protecting IP and information that contributes to industry competitiveness.”

Noting Google’s recognition of the importance of protecting its own IP, does Google also recognise the importance of creators protecting their IP and copyright from LLM developers?

Google recognises the importance of creators protecting their IP and copyrights and believes its approach to the development and deployment of its LLMs and generative AIs is consistent with those rights.

9. At the hearing Ms Longcroft justified the theft of copyrighted works by stating: “If we were to exclude works that are still under copyright, works and information, data, blog posts, and other information that forms part of that enormous corpus of data, that would mean that data relating to modern events or cultural or social issues such as LGBTQI rights, for example, would be excluded from those datasets. It is predictable that the models would then show bias or have gaps or ignorance about those interests and about that large and important part of our society. We train our models on that large corpus of publicly available data in order to ensure that they are providing the most socially beneficial uses in their outputs.”

Why should Governments permit large corporations to ignore or break laws because they have unilaterally made a judgement that it is “socially beneficial”, particularly where they coincidentally have a commercial interest in breaking those laws?

Google’s development and deployment of its LLMs and Generative AI models is conducted consistently with applicable laws.

10. Would Google be open to third parties being permitted to steal Google IP if those committing the theft claim it was socially beneficial to do so? If not, is this a double standard?

Please see our response to the prior question.

11. At the hearing, Microsoft and Amazon advised that they provide indemnification commitments to users of their AI products in the event they are sued for copyright infringement. Does Google offer a similar guarantee, and if yes what is the scope of that guarantee, and why has Google decided to (or not to) provide that guarantee?

Yes. Google has offered its Cloud customers a copyright indemnification. This was announced on October 12, 2023 in a [blog post](#) entitled “Shared fate: Protecting customers with generative AI indemnification”.

[Ends]