

CCIA Comments on Korea Copyright Commission Surveys¹ on Copyright and AI

I. Measures to secure legal use rights for AI training materials

Since AI training uses a large amount of copyrighted works, there are ongoing differences of opinion between copyright holders, AI developers, and service providers regarding the use of copyrighted works.

1) Can using copyrighted works for AI training without the permission of the copyright holder be considered fair use under current law? Why do you think so?

Yes, the use of copyrighted works for training AI models without permission can be considered fair use under current law. The processes for ingesting and training represent highly transformative fair uses. Ingesting expressive works serves the purpose of extracting unprotectable elements such as facts, patterns, and trends, and is not done for the purpose of copying or commercializing expression. Additionally, ingested data is often converted into usable formats and temporarily reproduced to protect against loss of data – such processes are not viewable or consumable by the public and therefore do not function as market substitutes for copies of the ingested works. With respect to large-scale models, while a piece of data may assist a model to better understand the meaning and inter-relationship of words, in general no individual piece of content has a particular influence over the model as a whole.

Leading U.S. copyright scholars have concluded, based on an examination of U.S. case law, that the ingestion of works for training purposes is the U.S. Copyright Act’s fair use provision (17 U.S.C. § 107), which is similar to Article 35-5 of the Korean Copyright Act.² Likewise, the Ministry of Justice in Israel has issued an opinion letter stating that ingestion of works for machine learning is permissible under the Israeli fair use provision, which also is similar to Article 35-5 and 17 U.S.C. § 107.³ While several dozen copyright infringement actions have been filed against AI developers in the United States, these cases are still in their early stages, and there is every reason to believe that AI developers will prevail in their assertion of fair use.

¹ Questions from the Surveys are in bold, CCIA answers are below each survey question. The surveys in question are AI Surveys 1-8, available here: <https://www.copyright.or.kr/customer-center/survey/list.do>.

² https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4654875;%20;https://texaslawreview.org/fair-learning/;%20;https://lawreview.law.ucdavis.edu/archives/53/2/copyright-and-progress-science-why-text-and-data-mining-lawful.

³ <https://project-disco.org/intellectual-property/011823-israel-ministry-of-justice-issues-opinion-supporting-the-use-of-copyrighted-works-for-machine-learning/>.

Training models on copyrighted works without permission should be considered fair use regardless of the commercial nature of the entity involved or the output. For the former, the use is highly transformative and has no impact on the market. For the latter, AI systems will largely have substantial non-infringing uses, and should therefore receive fair use protection. If the output is infringing, the user who prompted the infringing output should be liable for copyright infringement.

The Commission's clear declaration that the use of copyrighted material in training and deploying AI models is fair use would be welcome, and would benefit Korean innovators and citizens. The Commission could also consider developing an additional AI-specific exception to provide further regulatory certainty for system developers. This approach has been followed in Singapore, which has both a fair use provision and a specific provision targeted at text and data mining.⁴

2) Do you think that compensation should be paid when using copyrighted works for AI training, unless it is fair use?

Using copyrighted works for AI training should be considered fair use, so this question is largely resolved by the application of fair use explained above. Therefore, while compensation may be paid on a voluntary basis, as some AI developers have opted to do, it should not be a requirement under fair use of copyrighted works for AI training. The fact that licensing has emerged only in certain, specific circumstances may present evidence that the market is working, such that licensing is being used only in scenarios where there is particular value or where the license is primarily for non-copyright reasons such as paying for ease of access. Additionally, if an AI system is designed and intended to reproduce copyrightable expression in its output, that could create a basis for a claim for compensation.

3) If you think monetary compensation is necessary, please explain why and how it should be provided.

4) If you think that non-monetary compensation is necessary, please explain why and what kind of non-monetary compensation there might be.

5) What are some efficient ways to obtain permission from copyright holders when using large-scale data for AI training?

⁴ <https://project-disco.org/intellectual-property/021021-draft-singapore-copyright-bill-proposes-significant-innovations/>.

In general, there is no efficient, enforceable way to require consent from copyright holders when using large-scale data from open, publicly available sources, given the huge volume of works involved.

Direct voluntary licensing would be infeasible in almost all creative sectors, at least in combination with an opt-in system. Given the large volume of work produced on a daily basis, a licensing process would be unlikely to keep pace. In addition, a significant amount of material used to train AI models lacks any identifiable author for obtaining permission, and of the portion that has an identifiable author, contacting them might be difficult. Further, the author is not necessarily the rightsholder, meaning that even after an author is identified it may be difficult to contact the party with the right to license the copyrighted work.

Moreover, implementing a required consent approach would distort the market. Only AI developers with sufficient resources could license the necessary large-scale datasets from the most significant (and often most litigious) copyright owners, excluding smaller AI developers as well as smaller creators and resulting in lower quality, less diverse datasets. Meanwhile, the higher costs of obtaining permission for a smaller share of licensed materials would result in less capable and potentially more biased AI models. Mandating licensing agreements for generative AI would lead to inferior technologies, fewer competitors in the marketplace and hindered innovation generally.

Nevertheless, some AI companies are developing mechanisms to allow copyright holders effective choice in whether to allow their content for training, e.g., by enabling technical mechanisms for opting out using the well-understood and widely used robots.txt protocol. Any technical mechanisms must be both an industry standard and machine-readable, as robots.txt is, in order to enable innovation.

II. Institutional improvement plan related to disclosure of training data

Copyright holders believe that it is necessary to disclose a list of AI training data so that it is known which copyrighted works were used for AI training, but AI developers and service providers believe that it is difficult to disclose a list of AI training data due to reasons such as trade secrets or excessive administrative costs. European AI law stipulates that ‘a sufficiently detailed summary of the content used to train general-purpose AI models must be prepared and made publicly available.’

1) Do you think AI developers or service providers have an obligation to collect and keep records of materials used for training?

No, AI developers and service providers should not be generally required to collect or keep records of materials used for training under copyright law. First, there are no such obligations for individuals to keep track of the materials they read in order to learn. Second, such a requirement would be costly and technically difficult, especially given the extremely large

amount of materials involved, and would likely undermine research and innovation. Third, maintaining such records would provide little benefit given the application of fair use discussed above. Finally, as discussed above, many AI companies are adopting means to enable website owners from opting out of having their materials crawled for training, e.g., via the robots.txt protocol.

If the ingestion of a work is a fair use, what would be the point of maintaining a record of what works were ingested? If AI developers or service providers have the above storage obligations, do you think they should be able to view them upon the copyright holder's request?

See above.

2) Do you think AI developers or service providers have an obligation to disclose records of the materials used for training?

See above. Additionally, such obligations could implicate highly commercially sensitive information, including trade secrets or other confidential information.

3) If there is a disclosure obligation, do you think it is necessary to make the records of training data publicly available, or do you think it is sufficient to make them publicly available only to the copyright holder?

See above. Given the technical infeasibility of retaining such data in the first place, there is no meaningful distinction between requiring public disclosure versus disclosure to the copyright holder.

4) How much information do you think a “sufficiently detailed summary” should include? (e.g. links to sites studied, names of authors and works, dataset information such as “Common Crawl”, “Books3”)

III. Alternatives to pre-authorization of AI training works

According to the Copyright Act, prior permission must be obtained to use a work, except in special circumstances. AI training uses a large amount of data, and it is sometimes difficult to find the rights holder for some works. Therefore, it is difficult to obtain prior permission for use of all works for AI training. Therefore, there is an opinion that alternatives are needed for utilizing training data.

In Europe, text data mining for AI training is permitted, but if the copyright holder explicitly states that he/she does not allow the collection of his/her data (Opt-out), information collection is prohibited.

1) Some countries, such as Europe and Japan, have introduced regulations (hereinafter referred to as TDM exemption regulations) that exempt users from having to obtain prior permission to use copyrighted works through text data

mining if certain conditions are met. Do you think that this TDM exemption regulation should be introduced in Korea as well? What are the reasons for and against it?

Korea should adopt a TDM exemption similar to those already enacted by the EU, Japan, and Singapore. Such an exception would help spur innovation by providing greater certainty for AI developers in Korea. The TDM exception should be available to commercial and non-commercial entities alike, as much of the innovation in the AI ecosystem is being fueled by commercial entities. In addition, it is recommended that Korea adopt a TDM exception that does not include conditions that significantly restrict the scope and pro-innovation impact of the exception. For example, as noted above, AI developers are already adopting means to enable website owners from opting out of having their materials crawled for training, e.g., via the robots.txt protocol, making a specific condition with respect to opt outs unnecessary. This approach will still allow rightsholders to proactively opt out using machine-readable, industry-standard means. While some stakeholders oppose such exemptions on the basis of enacting more prescriptive measures, more stringent measures are often technically infeasible and would impose harms on rightsholders, AI developers, and the general public, as detailed in the below answers.

2) Looking at overseas cases, there are cases where the TDM exemption provision has certain conditions, such as (1) allowing only for non-profit and research purposes, or (2) only using works that the copyright holder does not object to. If the provision were introduced in Korea, what type of TDM exemption provision do you think would be necessary? Why?

See above.

3) If you are against the introduction of the TDM immunity provision, what are some alternatives to address this situation?

4) Do you think that in our country, copyright holders need to be able to express their opposition (opt-out) to the use of their copyrighted work for AI training?

Because the use of AI training for copyright work should be considered fair use, no affirmative consent should be required. However, copyright owners who wish to may have effective means of opting out of allowing their works as training materials. For example, some may be able to put their content behind technological protection measures such as paywalls. Copyright holders making content available on the web are able to use the widely-used robots.txt exclusion protocol to prevent the work posted to their websites from being crawled by specific AI bots.

The idea of an opt-in regime would, as stated above, block AI development and produce market distortion. Training data created in less common languages or from various subcultures is far less likely to be organized, and the appropriate entity to contact for permission may even be impossible to determine, while SMEs without the resources to identify and negotiate with rightsholders could be excluded from AI development.

5) What do you think would be an appropriate way for a copyright owner to express their objection to their work being used for AI training? Explain why

a) Method of controlling AI from crawling copyrighted works on a website (robot.txt, robot exclusion technology)

Any standard developed could be made voluntary for developers' compliance as a method of respecting reservations of rights.

b) If the author requests an opt-out after AI has trained on a copyrighted work, a method of returning to a state in which the work has not been learned (unlearning*)

c) An act of excluding data for which training rejection has been requested from among the data already learned by AI

d) Other methods

As mentioned above, the most appropriate method would be through controlling AI from crawling copyrighted works on a website through robots.txt or robot exclusion technology, for example. That protocol allows for granularity that would permit search engine bots but exclude other bots, or would permit a bot to ingest data from a site for some uses but not others. Several companies have recently announced extensions that will allow website publishers to allow search bots but exclude AI training bots.

The proposal to allow authors to retroactively request to have their copyrighted work “unlearned” would, at least at present, be technically impossible without fully retraining a model, as the specific impact on inferences derived from a particular piece of training material is not retained, to the extent it even exists in the first place. Given the large expense of retraining a model, including significant energy consumption, there is no economically feasible way to ‘unlearn’ inferences from a particular piece of training data. (Moreover, new releases of models typically are retrained from the ground up. Thus, if a copyright owner “opts out” after release 1.0 of a model is trained, while release 1.0 can’t be retrained, release 2.0 will not be trained on their work.)

This area of technology continues to rapidly develop, and it is possible that unlearning might become economically feasible in the future. At the same time, there is no guarantee that this circumstance will come to pass. Accordingly, the Copyright Commission should conduct its

analysis and make its recommendations based on the assumption that unlearning is not feasible.

- 6) What remedies are available to the author if the developer refuses to exclude data after the AI has trained on it? Is it enough to punish those who use the work without permission through existing copyright laws? Or should new regulations be created to regulate the act of refusing to exclude data after the AI has trained on it?**

While industry would not intentionally ignore such an objection, this issue has not been addressed by the courts in a manner that provides clear jurisprudence.

IV. Copyright infringement of output

Under copyright law, in order to constitute copyright infringement, it is judged based on whether ‘the original work is referenced (reliance on the original work)’ and ‘whether it is actually similar (substantial similarity)’.

- 1) If it is determined that an AI-created output infringes on the copyright of an original work created by humans, should it be judged according to existing standards as is now the case? Or should a different standard suitable for AI be used?**

Without commenting on the specifics of Korean law and jurisprudence in this space, CCIA suggests that any infringement analysis should be considered equal between humans and AI systems. Infringement, by any party, should be treated the same.

- 2) If AI output can infringe on the copyright of human-created works, who do you think is responsible for the infringement? Explain why**
 - a) AI developer (within the development company)**
 - b) A person who trained AI**
 - c) A person who used AI through prompt input**
 - d) A person who copied, distributed, published, or did similar acts with AI output**
 - e) Other**

As a general matter, any liability should be placed on the end-user who requests and publishes a copyright-infringing work, not the model, system, or developers that trained or deployed an AI system. Similar to other technologies that may implicate IP, such as video recording devices (i.e., VCRs) or computers used to replicate content without permission, the output of any use of AI systems should be attributable to the user whose volition sets into motion the actions

alleged to be infringing, not the provider of the system being abused by that end user, particularly if the end user is violating the terms of use of the provider.

V. Whether to display AI output

With the advent of generative AI, it has become realistically difficult to distinguish AI output from human works, and there is discussion about the need for marking to distinguish between works protected by copyright law and AI output that is not protected.

1) Do you think that labeling should be mandatory for AI-generated images, videos, texts, music, and other AI-generated outputs? Or do you think that it is not necessary? Why?

Practices are quickly developing in the industry, and the government should foster continuing innovation in the development of technologies and industry standards before rushing to impose mandatory obligations. As industry practices develop, it would bring helpful information on whether labeling is helpful for consumers and reasonable for companies, and, if so, what labeling has worked best.

A blanket requirement to label all AI-generated outputs is not feasible and likely unnecessary for consumers, particularly for text. For example, CCIA believes it clearly would not be appropriate to require the labeling of the output of a translation AI. Moreover, in general, norms should be allowed to develop within different fields concerning what labeling is appropriate, effective, and robust. Thus, there typically would be no need to label an AI-generated response to a basic consumer inquiry in a commercial interaction, but a student likely should indicate AI assistance in coursework.

If the government seeks to go down this path, CCIA recommends including a narrower scoping for verification tools. Best practices globally focus on providing consumers with the ways of recognizing or proactively identifying AI-generated content and content that has been materially altered by generative AI. Because technologies for embedding watermarks in text are still developing, however, any mandated labeling scheme should be narrowed to particular, urgent use cases.

If possible, and regardless of whether the digitally manipulated or created content is text, audio, photo, or video content, any labeling requirement should distinguish between fully AI-generated content and content that is materially changed in character and meaning through AI-powered tools, but such a requirement should not be imposed for immaterial or slight alterations of content that may incidentally be facilitated by AI-powered tools.

Additionally, the labeling of output should not be mandatory if an application is itself clearly identified as a generative AI app, or this identification could be included in metadata.

2) If labeling of AI output were required by law, what would be required to be displayed? (e.g., AI program used, human-worked portion, number of prompts, degree of collaboration between humans and AI (human 00%, etc.)

CCIA strongly advises that the government not adopt the proposed framework of “degree of collaboration between humans and AI” with a percentage affixed. Requiring users or developers to quantify how much of an output is generated through human action vs. AI systems would be very difficult in many circumstances, particularly as the use cases for AI are still nascent and the nature of collaborative outputs will still develop over time.

As such, focusing on affixing a less specific label, such as a disclosure that “an AI program was used in generating this output,” would be preferable both for deployers and users of AI and consumers viewing the content. Inclusion of poorly-calculated percentages of human vs. AI input would likely lack long-term accuracy as well as instilling confusion in users.

Additionally, if there were a requirement to include a label, any such obligation should include the option to add the label in metadata of an AI-generated output.

3) If labeling of AI outputs is mandated by law, who should be the labeling entity? (Example: AI model developers, AI service providers, AI users)

The responsibilities of labeling AI outputs should be shared between AI providers and deployers, including users. AI model providers should offer users the technical tools needed to label AI-generated content, while users should the obligation to actually label AI-generated content should fall on consumers.

4) Do you think that labeling is necessary for all content creation using AI, or is it better to limit it to certain cases? (e.g., in the case of deepfake technology, in the case of false impressions, in the case of content directly related to human life, etc.)

As noted previously, it would be prudent to allow industry practices to evolve here to review what works best for consumers and businesses before mandating all content created with AI be labeled as such. If a standard is imposed too early, but as the technology and its use cases evolve the standard is discovered to be flawed or difficult for businesses to comply and/or consumers to use, this would be more harmful to the goal of maximizing transparency.

If an AI disclosure requirement were in place for all AI-generated content, it would likely quickly become overwhelming for consumers when multiple AI tools are used together to create a given piece of material or where tools incorporate AI technology but are not fully generative. For example, tools like code completion in development environments or content-aware fill in photo editors may utilize AI technology, but the ultimate output is not what consumers would understand to be “AI-generated.” For consumers, this sort of disclosure would be confusing and would lead most users to instead skip through the disclosure or ignore it.

Further, as previously stated, one particular case where the labeling of output should not be mandatory is if an application is itself clearly identified as a generative AI app. If a requirement is adopted anyway, this identification should be permitted to be included in metadata.

- 5) Do you think there needs to be a separate sanction for violations of the labeling obligation? If so, what level do you think would be appropriate?**
- a) Administrative measures such as fines, surcharges, etc.**
 - b) Criminal sanctions such as fines**
 - c) A method of self-regulation by establishing declarative and cautionary provisions without separate sanctions (e.g., a system that allows YouTube to restrict monetization of creators who violate labeling regulations or delete related content, etc.)**
 - d) Other**

As noted above, labeling requirements at this stage are premature.

- 6) There are two ways to introduce a display system for outputs: one that can be seen by anyone with the naked eye, and one that can be seen through technical processing to prevent damage to the aesthetics of images or videos. If a display system is introduced, what level do you think it should be introduced at?**

CCIA does not advocate for any specific form of labeling, and instead suggests that the issue raised in this question reflects precisely why the government should ensure that industry is able to implement rules that are consumer-friendly and are seamlessly integrated into their services when appropriate. For some cases, a label visible to the naked eye may make sense to the provider or user, while in others, embedding the label in metadata may be preferred. It is also possible that more than one of these techniques may be used simultaneously. Requiring one form of disclosure or another for certain instances would undermine the burgeoning development of AI content creation, and could result in users either ignoring the disclosure or certain uses of AI failing to realize their potential, due to visible labels that undermine the artistic aesthetics of an AI-generated piece.

- a) A method of displaying something visually so that people can see it;**
- b) A method that is difficult to distinguish with the naked eye but can be confirmed through technical processing**
- c) Other methods (e.g., in principle, make it possible to confirm with the naked eye, but apply a technical processing method depending on the situation).**

VI. Copyright registration of output

According to the interpretation of the current copyright law, the work itself created by AI, not humans, is not recognized as a work of authorship. However, if a human's creative contribution is added to the AI work through modification, addition, editing, arrangement, etc., the work is recognized as a work of authorship for that part, and the work can be registered. The effect of registration is limited to the part to which the human made a creative contribution.

1) Are the copyright registration standards for AI outputs like the above appropriate? Inappropriate? Why?

This would be sufficient. For example, under current U.S. Copyright Office guidelines, humans who use AI to create a work “may claim copyright protection for their own contributions to that work,” excluding any AI-generated content that is more than de minimis. This provides certain protections to an end user, so long as the human exercised adequate creative control over the work’s expression, and “actually formed” the traditional elements of authorship. Such a threshold is appropriate for granting copyright protection to AI-generated works.

Adequate contribution could include either a human author significantly changing the AI’s output into a final work, or by a human author exerting sufficient control over the output of a generative AI. However, so-called “prompt engineering”—requesting an AI system to produce a certain product through requests—should not alone be deemed sufficient for copyright protection on the final work.

2) If the current system is adequate, what are some specific ways to prove the additional creative contribution of humans?

3) If the current system is inadequate, what factors do you think should be considered when registering copyright for AI output?

4) If there are any areas in which improvements are needed in the current copyright registration system, copyright law, or registration manual in relation to artificial intelligence, please feel free to describe them.

VII. Protection of AI outputs without human creative contribution

There are various conflicting opinions on the protection of AI outputs without human creative contribution. Specifically, there are opinions that AI outputs should be protected

like copyrighted works, opinions that they should be protected to a lesser extent than copyrighted works, and opinions that there is no need for protection at all.

1) Do you think that AI outputs without any creative human contribution should be protected by law? Or should they be freely available? Why?

A work produced by an AI algorithm or process, without the involvement of a person contributing to the output, should not qualify as a work of authorship. If an AI algorithm or process creates a work that lacks expression from a person in the resulting output, this should not be registered as a copyrighted work.

Withholding copyright protection from a work that is generated through an AI process for which there was no tangible or expressive contribution by a person is reasonable from a policy perspective. The AI algorithm, and the computer that powers it, does not need the economic incentive that comes with a copyright in order to produce works. Indeed, AI is capable of quickly generating a huge number of outputs in a short time period. Therefore, granting copyright to AI-generated output could quickly create a broad swath of legal issues that bring litigation and uncertainty to AI developers and users.

2) Current law does not protect AI outputs that do not involve human intervention. What problems could arise if such AI outputs are not protected?

3) Do you think that if AI output without human intervention is to be protected legally, it should be at the same level as human creations? If so, what level of protection do you think is needed?

4) If AI output is to be protected, who do you think should have the rights to it? Explain why

- a) **AI developer (within the development company)**
- b) **A person who trained AI**
- c) **A person who created AI output through prompt input**
- d) **A person who copied, distributed, published, or performed similar acts on AI output**
- e) **Other**

The legal battles over who should receive compensation and copyrightability from works produced solely by AI are precisely why these works should not receive copyright.

VIII. Human creative contribution and protection of AI output

Under current law, outputs created through generative AI cannot be recognized as works. However, if a user modifies, adds, or adds to an output created by AI, the creative contribution can be recognized as a work of authorship.

Under copyright law, an author refers to a person who has made a creative contribution to a work, and when determining the work of authorship, whether or not there has been a creative contribution is one of the main factors in the judgment. In relation to this, there is an opinion that the act of inputting prompts in the process of creating an AI output should be considered a creative contribution.

- 1) If human creative contributions are acknowledged in AI output, who do you think should be the owner of the rights to the AI output? Explain why**
 - a) AI developer (within the development company)**
 - b) The person who trained the AI**
 - c) The person who created the AI output through prompt input**

The person who leverages an AI model or system to create a new work—again, with the caveat that mere “prompt engineering” should not be considered a sufficient exercise of creative control—should be acknowledged, if any one individual is to be acknowledged. Other persons—such as AI developers—can receive acknowledgement at other stages of the AI life cycle (copyright over code, compensation for development, etc.), but when focusing on an AI output specifically, the individual that created that output with the help of AI should be the one acknowledged.

- d) The person who copied, distributed, published, or did similar acts to the AI output**
 - e) Other**
- 2) Can the act of a user inputting a prompt into a generative AI be considered a creative contribution to the output produced by the AI, or is it simply an idea?**

As detailed earlier, “prompt engineering”—requesting an AI system to produce a certain product through requests—should not alone be deemed sufficient as a “creative contribution,” nor should it be deemed adequate for copyright protection on the final work.

- 3) If we consider the user's prompt input behavior as a creative contribution, to what extent can it be considered a creative contribution? (e.g., the amount and content of instructions/inputs, the number of times they were created, etc.)**



4) What are some ways to verify user prompt input behavior that is leveraged in AI output?