



November 17, 2023

## BSA COMMENTS ON ARTIFICIAL INTELLIGENCE AND COPYRIGHT GUIDELINES

### Submitted Electronically to the Korea Copyright Commission

BSA | The Software Alliance (**BSA**)<sup>1</sup> welcomes the opportunity to provide our comments on the AI and Copyright Guidelines (**Guidelines**) prepared by the Korea Copyright Commission (**KCC**). BSA is the leading advocate for the global software industry before governments and in the international marketplace, and we have long supported effective copyright protection. Our members are at the forefront of providing AI-enabled products and services, as well as tools used by others in the development of AI systems and applications.

AI can benefit all industry sectors — including content creators — consistent with copyright law. AI systems learn from the computational analysis of large amounts of data that can be used to train a model. The model will then make predictions when presented with a query that includes new data. Some generative AI models can take this a step further and generate new data, such as text or a picture or even new code. Regardless of what form of AI is being used, respect for copyrighted works — as well as personal data privacy and trade secret information — should be a cornerstone of responsible AI development. This critical advancement in technology need not be at the expense of artists and rightsholders.

### Summary of BSA's comments

BSA highlights the following for the KCC's consideration:

1. **First**, for training AI models, computational analysis typically involves turning data into tokens that are statistically correlated with other tokenized data. Some of that data may be part of a copyrighted work, but the use of the data is not related to the expressive content of a work. A book may be used to learn language skills, which are then used for improved database management practices. This use falls within Articles 35(2) and (5) of Korea's Copyright Act (**Copyright Act**).<sup>2</sup> Furthermore, we understand there are voluntary industry conversations

---

<sup>1</sup> BSA's members include: Adobe, Alteryx, Altium, Amazon Web Services, Asana, Atlassian, Autodesk, Bentley Systems, Box, Cisco, Cloudflare, CNC/Mastercam, Dassault, Databricks, DocuSign, Dropbox, Elastic, Graphisoft, IBM, Informatica, Juniper Networks, Kyndryl, MathWorks, Microsoft, Nikon, Okta, Oracle, Palo Alto Networks, Prokon, PTC, Rockwell, Rubrik, Salesforce, SAP, ServiceNow, Shopify Inc., Siemens Industry Software Inc., Splunk, Trend Micro, Trimble Solutions Corporation, TriNet, Twilio, Unity Technologies, Inc., Workday, Zendesk, and Zoom Video Communications, Inc.

<sup>2</sup> Copyright Act, amended by Act No. 19410, May 16, 2023, translated in Korea Legislation Research Institute online database: [https://elaw.klri.re.kr/eng\\_service/lawView.do?hseq=62935&lang=ENG](https://elaw.klri.re.kr/eng_service/lawView.do?hseq=62935&lang=ENG) (**Copyright Act**).

Article 35(2) (Temporary Reproduction in Course of Using Works) reads: *Where a person uses works, etc. on a computer, he or she may temporarily reproduce such works, etc. in that computer to the extent deemed necessary for the purpose of smooth and efficient information processing: Provided, That this shall not apply where the use of such works, etc. infringes on copyright.*

Article 35(5) (Fair Use of Works) reads: *(1) Except as provided in Articles 23 through 35-4 and 101-3 through 101-5, where a person does not unreasonably undermine an author's legitimate interest without conflicting with the normal exploitation of*

around developing and using automated tools to indicate when a rights-owner does not want content used for training purposes, similar to the current “do not crawl” tools. We encourage further conversations to determine whether a consensus standard is possible.

2. **Second**, for the copyrightability of works generated using AI systems, the analytical touchstone should be whether human creativity was responsible for the work, regardless of what instrument or technology was used to aid its expression. This is consistent with established practice where software tools are used by artists to create works that are regularly deemed copyrightable, even though their creation was aided by technology. Generative AI should bolster creativity, just as other software applications have long been an important tool of artists and storytellers. Generative AI is used, for example, in word processing by authors and photo enhancements by visual artists; it is used to create special effects in audio-visual works and arranging music for sound recordings; in software development, it is used to assist in generating software code based on the programmer’s instructions. When generative AI is used to enhance human creativity, the resulting work should be protected by copyright. In instances where AI-generated works do not contain creative elements, but are combined with human-authored works, it should not render the entire combined work unprotectable. Instead, the otherwise unprotectable AI-generated portions should be disclaimed, but protectable portions of the combined work should be copyrightable.
3. **Third**, as regards the risk that AI systems will be used to create infringing works, the copyright laws are sufficiently flexible to determine whether copyright infringement liability exists. Stated differently, whether the output is infringing is typically not related to the technology used. Where an output is infringing, there should clearly be liability regardless of whether AI is used in producing the infringing copy.

## AI in the copyright context

AI machine learning encompasses a vast array of technologies developed or deployed for use in a variety of different industries and applications. Machine learning depends upon the computational analysis of training data to identify correlations, patterns, or other metadata to develop a model that can make predictions or recommendations based on future data inputs.<sup>3</sup> More recently, generative AI models have emerged which are able to generate new text, image, or sound.

### Applications of AI

To explain the insights, predictions, and other outputs derived from computational analysis in the machine learning context, we provide below a few widely recognized examples:

- Automated flight management and air traffic control based on computational analysis of meteorological conditions, real-time fuel consumption, aircraft operational data, nearby air traffic conditions, airport congestion, and numerous other data elements.<sup>4</sup>

---

*works, he or she is entitled to use such works. (2) In determining whether an act of using works falls under paragraph (1), the following matters shall be considered: 1) Purposes and characteristics of use; 2) Types and purposes of works; 3) Amount and substantiality of portion used in relation to the whole works; 4) Effect of the use of works on the existing or potential market for the works or current or potential value thereof.*

<sup>3</sup> BSA | The Software Alliance, *Confronting Bias: BSA’s Framework to Build Trust in AI* (2022) (**BSA AI Bias Framework**), at <https://ai.bsa.org/confronting-bias-bsas-framework-to-build-trust-in-ai>; See also, National Institute of Standards and Technology, *NIST Risk Management Framework* (2023), at <https://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.100-1.pdf> (describing an AI system “as an engineered or machine-based system that can, for a given set of objectives, generate outputs such as predictions, recommendations, or decisions influencing real or virtual environments”).

<sup>4</sup> M. Durgut, *Artificial Intelligence and Air Traffic Control*, Aviationfile.com website (Jan. 2023), at [https://www.aviationfile.com/artificial-intelligence-and-air-traffic-control/#:~:text=One%20of%20the%20primary%20applications%20is%20to%20help.make%20informed%20decisions%20on%20routing%20and%20scheduling%20flights](https://www.aviationfile.com/artificial-intelligence-and-air-traffic-control/#:~:text=One%20of%20the%20primary%20applications%20is%20to%20help.make%20informed%20decisions%20on%20routing%20and%20scheduling%20flights;);

- Identification of chemical and cellular anomalies for early diagnosis, prevention, and treatment in the fields of oncology, autoimmune disorders, and Parkinsons and Alzheimer's disease.<sup>5</sup>
- By integrating generative AI into security operations, organizations can effectively identify and address security anomalies, as well as detect and mitigate potential threats.
- Predictive climate modeling based on computational analysis of satellite data, weather station data, topographical information, and various IoT and sensor data.<sup>6</sup>
- Improved carbon tracking and mitigation based on computational analysis of transportation logs, meter readings, fuel purchase records, atmospheric pollution tracking, and visual monitoring of power plants and other facilities, and other data sources.<sup>7</sup>
- Computational analysis to map vulnerable seaside areas to produce cyclone risk maps and guide investment plans for cyclone shelters, schools, health facilities, and other infrastructure for disaster planning and survivability.<sup>8</sup>
- Predictive typing and other office productivity solutions (e.g., an "auto-complete" function that suggests the letters "...cerely yours" after the typist inputs the letters, "sin"),<sup>9</sup> or the creation of sound effects or special effects to assist creators and film producers in developing new artworks.<sup>10</sup>
- Generative AI can also enable conversational interfaces to enable people to work with AI more easily and flexibly across a broad array of domains and applications, such as the applications referred to above. Conversational interfaces can enable people to query and reason over data and outputs of AI systems.

While the use cases are diverse, the elements of training each are very similar as further discussed below. The purpose of training is to enable the model to learn the unprotectable elements of copyright works they are trained on.

### The AI development life cycle

The AI development life cycle typically includes the following steps:

- Project Conception: First, the AI development team will formulate the "problem" that a system is intended to address and map the structure and target variables that the system is intended to predict. For models trained for a particular task, sometimes referred to as narrow AI applications, this may be a fitness app that analyzes a consumer's heart rate to monitor for irregularities that might predict whether that person is at risk of a stroke or heart disease (i.e.,

---

<sup>5</sup> Hunter et al., *The Role of Artificial Intelligence in Early Cancer Diagnosis*, 14(6) *Cancers* 1524 (2022), at <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8946688/>; Stafford et al., *A systematic review of the applications of artificial intelligence and machine learning in autoimmune diseases*, 3 *NPJ - Digital Medicine* 30 (2020), at <https://www.nature.com/articles/s41746-020-0229-3>; Diogo et al., *Early diagnosis of Alzheimer's disease using machine learning*, 14 *Alzheimers Research and Theory* 107 (2022), at <https://alzres.biomedcentral.com/articles/10.1186/s13195-022-01047-y>.

<sup>6</sup> Schneider et al., *Harnessing AI and computing to advance climate modelling and prediction*, 13 *Nature Climate Change* 887 (2023), at <https://www.nature.com/articles/s41558-023-01769-3>; World Economic Forum, *The role of machine learning in helping to save the planet* (2021), at <https://www.weforum.org/agenda/2021/08/how-is-machine-learning-helping-us-to-create-more-sophisticated-climate-change-models/>.

<sup>7</sup> Global Data Alliance, *Cross-Border Data Transfers & Environmental Sustainability* (2023) (internal citations omitted), at <https://globaldataalliance.org/wp-content/uploads/2023/04/04192023qdacbdtsustainability.pdf>.

<sup>8</sup> Global Data Alliance, *Cross-Border Data Transfers & Environmental Sustainability* (2023).

<sup>9</sup> S. Ashraf, *Desmystifying Autocomplete*, *Towards Data Science* (2020), at <https://towardsdatascience.com/index-48563e4c1572>; A. Wickramarachchi, *Machine Learning: Word Embedding and Predicting*, *Towards Data Science* (2020), at <https://towardsdatascience.com/machine-learning-word-embedding-and-predicting-e603254e4d7b>.

<sup>10</sup> See e.g., D. Nelson, *AI Researchers Design Program To Generate Sound Effects For Movies and Other Media*, *Unite AI Website* (2022), at <https://www.unite.ai/ai-researchers-design-program-to-generate-sound-effects-for-movies-and-other-media/>.

the target variable). Foundation models will however be developed to power a broad array of applications.

- **Raw Data Identification:** Second, the AI development team may identify a relevant universe of “raw data” that will be subsequently transformed and structured. Data sources are as diverse as the potential applications of machine learning AI and may include everything from machine-to-machine data (e.g., satellite transmission data) and international trade statistics to published materials, blog posts, website comments, and chat room logs. “Raw data” is frequently “messy,” requiring significant work to transform the data into a usable form, as outlined below. The data to develop a narrow AI system for a specific application or fine-tuning a foundation model for a particular task will relate to the particular task. Foundation model training will require broader and more varied data. The scale of data needed to train large language models is vast and will rely on being able to analyze data on the Internet to achieve this scale.
- **Preparing the Data Set:** The AI development team may modify the “raw data” so that it can be understood semantically by the machine and used to train the model. During this process, the team will revise, clean, and normalize the data as necessary. Data typically is transformed semantically and structurally through “tokenization,” which involves breaking down a piece of text or data into smaller units (or “tokens”) for purposes of computational analysis. Additional processing may be necessary to improve the reliability, quality, and suitability of the data for analysis, helping to address quality challenges such as missing values, duplicates, outliers, and inconsistent formatting across the entire data set. Large scale AI models are typically trained using self-supervised methods, dispensing with the need to label the training data, and allowing the models to scale on a vast scale of training data typically collected from the Internet.
- **Model definition:** After input data has been suitably processed, the AI development team must establish the system’s underlying architecture. This includes identifying the variables (i.e., “features”) in the training data that the algorithm will evaluate as it looks for patterns and relationships as the basis of a rule for making future predictions. It also includes selecting the type of algorithmic model that will power the system (e.g., linear regression, logistic regression, deep neural network.)<sup>11</sup> Once the data is ready and the algorithm is selected, the team will train the system to produce a functional model that can make predictions about future data inputs.
- **Model Validation, Testing, and Revision:** After the model has been trained, the AI development team must validate it to determine if it is operating as intended and test it to demonstrate that the system’s outputs fall within expected parameters and do not contain unexpected errors or unintended bias. Based on the outcome of validation and testing, the team may need to revise and refine the model to mitigate these risks.

## The use of copyrighted works to train AI models

As discussed above, computational analysis is typically applied to a large training data set that may comprise millions or billions of tokenized data elements. Depending on how the model is trained, data accessible over the Internet may be collected as part of the raw data set that is transformed into the tokenized elements.

This raw data may include copyrighted works because a substantial portion of the content freely available on the Internet is potentially subject to copyright protection, which has a low threshold to establish “originality” and which provides that copyright arises automatically upon a work’s creation, even if it the work is not registered. Importantly, however, not all the material online is subject to

---

<sup>11</sup> BSA AI Bias Framework (2022).

copyright, in part because copyright protection does not extend to facts, ideas, or mathematical concepts.<sup>12</sup>

Computational analysis may involve two sets of reproductions that potentially implicate the Copyright Act: (1) reproductions necessary to create a set of “training data,” and (2) temporary reproductions that are incidental to the computational process of training the AI model. In each case, the reproductions are not visible or otherwise made available to the public. Instead, the reproductions are the necessary byproduct of a technical process that is aimed at identifying non-copyrightable information *about* the underlying data derived from the works — i.e., the correlations and patterns that inform the creation of the AI model and enable it to make predictions based on future data inputs. Such non-expressive reproductions are a) temporary reproductions, per Article 35(2) of the Copyright Act; and b) do not undermine an author’s legitimate interests that copyright is intended to protect, per Article 35(5) of the Copyright Act.

Furthermore, computational analysis does not involve the consumption of any copyrighted works for their expressive content. Rather, such analysis involves mathematical calculations of probabilities, correlations, trends, and other patterns across the entire tokenized data set. Such analysis seeks to understand only the mathematical patterns (e.g., the relationships of specific tokens in relation to other tokens) distributed across the entire data set. These mathematical patterns are themselves not expressive content protected by copyright law.

### Computational analysis of AI training data is fair use

The exceptions provided in Articles 35(2) and (5) of the Copyright Act ensure the use of copyrighted works for the purposes of analyzing large collections of information to identify patterns, correlations, and other metadata to develop AI models that makes predictions about future data inputs.

By way of illustration, this means that an AI developer seeking to create a natural language processing model — such as an AI-driven predictive typing model — can rely on publicly available text-based material to create the training database. In such a scenario, the AI developer would not be reproducing this text for its expressive purpose. Rather, the reproductions would be made solely for the purpose of extracting unprotected information about the correlations, patterns, and relationships among letters and words as they appear in thousands of phrases, figures of speech, similes, metaphors, grammatical patterns, and common linguistic formulations and expressions. Neither the letters, words, and phrases, nor the mathematical patterns among them across thousands or millions of writings, are copyright protectable subject matter.

Similarly, for an image generation model, the model developer will take a very large volume of images tagged with words — for instance, some number of cat photos tagged with “cat,” and dog photos tagged with “dog.” Over time, the model learns that certain patterns are characteristic of “cattiness”, and it will learn to recognize whether an image fed to it is a cat or not, whether the image is a cartoon, or a photograph, etc. Once again, the machine is not reproducing for expressive use any particular image of cats but is instead dissecting the images to understand what a cat is — a basic set of facts/statistical correlations rather than expressions.

Furthermore, the ultimate use of the computational analysis applied to the data set is a transformative use of the original content. Auto-complete functionality in predictive typing software comprises a new creation in the form of software code that is distinct both from the entire body AI training data and from any single work within that data set. Software code that can suggest the endings of commonly used phrases is many steps removed from any copyrighted works found in the underlying raw data.

Finally, such functionality simply does not compete with any copyrighted works in any manner that copyright is intended to protect. Accordingly, an AI developer can rely on the exceptions in Articles 35(2) and (5) to construct an AI training database using text, images, or other data to which she has lawful access.

---

<sup>12</sup> Article 9.2 of the TRIPS Agreement provides that copyright protection shall extend to expressions and not to ideas, procedures, methods of operation or mathematical concepts as such. In other words, copyright protection does not cover any information or ideas contained in a work; it only protects the original way in which such information or ideas have been expressed. Thus, everyone is free to use the information contained in a work, including for the purpose of creating new works.

## The copyrightability of material generated by AI

As for the copyrightability of works generated using AI systems, the analytical touchstone should be whether human creativity was responsible for the work regardless of what instrument or technology was used to aid its expression. Generative AI should bolster creativity, just as other software applications long have been an important tool of artists and storytellers. Generative AI is used, for example, in word processing for authors and photo enhancement for visual artists; it is used to create special effects in audio-visual works and arranging music for sound recordings; in software development, it is used to assist in generating software code based on the programmer's instructions. When generative AI is used to enhance human creativity, the resulting work should be protected by copyright. In instances where AI-generated works do not contain creative elements, but are combined with human-authored works, it should not render the entire combined work unprotectable. Instead, the otherwise unprotectable AI-generated portions should be disclaimed, but protectable portions of the combined work should be copyrightable.

As inquiries regarding copyrightability will turn on a close examination of the degree of human creativity, and with AI increasingly being used as a tool in all categories of creative works, a decision to limit copyrightability when AI is used would significantly chill adoption of AI solutions.

Works that emerge as outputs of AI systems and meet the human creativity requirement should continue to be eligible for copyright protection. In most cases, AI systems will function as tools used by human authors and creators to execute upon their creative vision. For instance, photographers will use AI-enabled tools to automate the tedious process of editing their images,<sup>13</sup> architects will use AI to augment their designs to enhance their energy efficiency,<sup>14</sup> and filmmakers will use AI to ensure that the movement of their animated characters appear more life-like.<sup>15</sup> In each of these cases the creative contribution of the human user makes it easy to conclude that the output would be copyrightable.

The use of generative AI should not change the analysis.

## Potential liability for infringing works generated by AI

As for the potential infringement of works generated using AI systems, current copyright laws are sufficiently flexible to determine whether copyright infringement exists. Whether the output is infringing should not be related to the technology used — where an output is infringing, there should be liability regardless of whether AI is used to produce the infringing copy. This is also consistent with established law, which focuses on whether a work infringes and is less concerned with the method used to create the work.

### Current copyright law protects copyright holders from infringement, including in cases arising from AI generated content

Copyright holders should have full and effective remedies when their rights are infringed. This principle applies equally to outputs generated using AI systems and outputs generated in other ways. Whenever such infringement is found, it is critical to compensate fully artists and creators for any damages caused.

In our view, existing copyright law should prove adequate to address questions of infringement. In most AI use cases, the output of an AI system will not implicate copyright at all. However, AI, like other technologies, could be used to create infringing material. In those situations, infringement liability would be premised upon proof that the allegedly infringing output is based on “actual copying” and is “substantially similar” to the copyrighted work, such that an “ordinary reasonable person would fail to differentiate between the two works” in view of the “qualitative[ly] and quantitative[ly] significant” similarity between the works.

In cases in which the copyright claims are “thin,” because there is only one way, or there are only a few ways, to represent specific facts, it would be necessary to prove “virtual identity” between the

---

<sup>13</sup> <https://theblog.adobe.com/adobes-general-counsel-makes-the-case-for-ai/>.

<sup>14</sup> <https://www.autodesk.com/redshift/machine-learning-in-architecture/>.

<sup>15</sup> <https://theblog.adobe.com/state-of-ai-in-animation/>.

output and copyrighted work, consistent with prevailing legal norms. Many users of generative AI will be small businesses experimenting with new tools and enforcement decisions should take account of this. Liability would also arise in other appropriate cases, such as those involving derivative works.

Plaintiffs may also seek to bring infringement actions against providers of AI-related services. Here too, the existing Copyright Act should prove adequate in evaluating and apportioning liability. If a plaintiff demonstrates that a direct infringement has occurred, courts will evaluate whether the service provider should be deemed liable for aiding or abetting copyright infringement.

### **AI models trained on sufficiently large data sets are less likely to produce infringing outputs**

While it is important not to conflate training data with the output of an AI system, it is worth noting that the more data available for training, the less likely the system will produce a copy or derivative of any particular input (in the absence of a user's efforts to do so). In any well-designed AI model trained on a sufficiently large data set, computational analysis should never (or only in the rarest of circumstances) produce outputs that are "substantially similar," let alone "virtually identical," to any specific copyrighted work.

Some AI developers and deployers are already taking steps to engage with artists and creators on how to support their work in a changing digital environment and taking steps to limit misuse of an AI system, such as limiting what prompts can be used.

### **Conclusion**

We hope that our comments will assist the KCC in its development of the Guidelines. As discussed above, the non-consumptive computational analysis of content in AI training data sets constitutes "fair use" under Korea's Copyright Act. BSA also supports multi-stakeholder efforts relating to AI training processes as well as efforts to minimize the risk of infringement. To the extent that infringement occurs, BSA strongly supports fully protecting content creators.

Thank you for the opportunity to share our views on these important issues, and please do not hesitate to contact me if you have any questions regarding this submission or if I can be of further assistance.

Sincerely,



Tham Shen Hong  
Manager, Policy – APAC