



25 October 2024

## BSA COMMENTS COPYRIGHT AND AI REFERENCE GROUP TRANSPARENCY PAPER

### Submitted Electronically to the Attorney-General's Department Copyright and AI Team

BSA | The Software Alliance (**BSA**) welcomes the opportunity to provide our response to the Attorney General's Department's Paper *Copyright and AI Reference Group: Transparency, September 2024* (**CAIRG Paper**).

BSA is the leading advocate for the global software industry. BSA members<sup>1</sup> create technology solutions that power other businesses, including cloud storage services, customer relationship management software, human resources management programs, identity management services, security solutions, and collaboration systems. Our members are global leaders in developing, tailoring, integrating, and deploying artificial intelligence (**AI**) systems and services, and the tools used by others in the development of AI systems and applications. As a result, they have unique insights into the technology's tremendous potential to spur digital transformation and the policies that can best support the responsible use of AI.

As we explain in detail in our submission<sup>2</sup> to the Department of Industry, Science and Resources (**DISR**) on its Proposals Paper for Introducing Mandatory Guardrails for AI in High-Risk Settings (**Proposals Paper**),<sup>3</sup> BSA supports encouraging transparency and information sharing across the AI value chain. However, any guidelines or requirements for information sharing should acknowledge that different entities possess distinct knowledge and abilities to share information. Furthermore, such guidelines or requirements must respect the need to protect confidential information, including trade secrets.

### Summary of Recommendations

Our comments in response to the CAIRG Paper focus on two overarching questions: 1) potential transparency obligations regarding the data and content used to train AI systems, and 2) potential

---

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

<sup>2</sup> Australia: BSA Comments on Mandatory Guardrails for AI in High-Risk Settings (BSA Comments), 03 October 2024, at: <https://www.bsa.org/policy-filings/australia-bsa-comments-on-mandatory-guardrails-for-ai-in-high-risk-settings>

<sup>3</sup> Safe and responsible AI in Australia: Proposals paper for introducing mandatory guardrails for AI in high-risk settings: (Proposals Paper), September 2024, at: [https://storage.googleapis.com/converlens-au-industry/industry/p/prj2f6f02ebfe6a8190c7bdc/page/proposals\\_paper\\_for\\_introducing\\_mandatory\\_guardrails\\_for\\_ai\\_in\\_high\\_risk\\_settings.pdf](https://storage.googleapis.com/converlens-au-industry/industry/p/prj2f6f02ebfe6a8190c7bdc/page/proposals_paper_for_introducing_mandatory_guardrails_for_ai_in_high_risk_settings.pdf)

transparency obligations related to AI outputs. Our responses and recommendations are summarized here and discussed further below.

- The Proposals Paper already addresses the importance of data disclosure and governance requirements, including transparency, for data used in AI training through Guardrail 3 (addressing data quality and provenance to protect AI systems).<sup>4</sup> Such transparency measures are appropriately focused on high-risk uses of AI, aligning with a risk-based approach that prioritizes safety and accountability where they matter most.

Transparency mandates for copyrighted content used in AI training should not be necessary and are not appropriate given that the use of on-line content, whether copyright protected or not, for AI training does not involve the consumption of works for their expressive content. Instead, AI training data is analysed mathematically to determine underlying patterns in the content, and such information is not subject to copyright protection. Copyright protection extends to creative expression, not the mathematical or statistical processing of that expression.<sup>5</sup>

Furthermore, a substantial portion of the information on the Internet, e.g. online blogs, reader commentary, chat room exchanges, and anonymous restaurant and hotel reviews, are potentially subject to copyright protections. As a result, serious questions of administrability persist for copyright-related transparency measures. Identifying the copyright status of every data point within a training dataset would be technically infeasible, particularly when considering the frequent use of anonymous or unattributed content.

For these reasons, **we urge the Australian Government to focus transparency measures on the high-risk AI context and avoid imposing such mandates for copyrighted content** before resolving administrability concerns.

- **We recommend the Australian Government amend the Copyright Act to adopt an explicit computational data analysis exception** to promote AI development and wider economic growth, facilitate research, allow Australia to maintain its competitiveness on the global stage and sustain its position as a hub for growth and investment, afford legal certainty to users, and strike an appropriate balance between copyright protection and reasonable use of copyright works.
- BSA supports the intent of the Proposal Paper's Guardrail 6 to inform end-users regarding AI-enabled decisions and interactions with AI and AI-generated content.<sup>6</sup> Specifically, **we support content authentication and provenance mechanisms**, such as the Content Authenticity Initiative's (CAI)<sup>7</sup> efforts to promote the open Coalition for Content Provenance and Authenticity (C2PA)<sup>8</sup> standard to help users identify AI-generated content.

Before discussing our responses to these questions, we provide a summary of how AI models are trained and how this pertains to copyright protections.

---

<sup>4</sup> Proposals Paper, pages 37-38

<sup>5</sup> TRIPS Agreement at: [https://www.wto.org/english/docs\\_e/legal\\_e/27-trips\\_01\\_e.htm](https://www.wto.org/english/docs_e/legal_e/27-trips_01_e.htm)  
Article 9(2): "Copyright protection shall extend to expressions and not to ideas, procedures, methods of operation or mathematical concepts as such."

<sup>6</sup> Proposals Paper, pages 39-40

<sup>7</sup> Content Authenticity Initiative at <https://contentauthenticity.org/>

<sup>8</sup> Coalition for Content Provenance and Authenticity at <https://c2pa.org/>

## Use of Copyright Protected Works in AI Training Data Sets

Large AI models are trained on massive amounts of unlabelled data through self-supervised learning using a vast corpus of training data and may later be customized for specialized tasks. In circumstances where the model is developed for a specific task, the AI development team additionally must identify a relevant universe of “raw data” that will subsequently be transformed and structured. Data sources are as diverse as the potential applications of the AI system 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. This raw data may include copyrighted works because a substantial portion of content on the Internet is potentially subject to copyright protection, as copyright protection typically applies upon a work’s creation.

It frequently requires significant work to transform this raw data into a usable form. During this process, the development team will revise, clean, and normalize the data as necessary. The data is typically transformed through “tokenization,” which involves breaking down text or other data into small units (or “tokens”) that are numerical vectors representing mathematical/statistical relationships within the training data for purposes of computational analysis. The tokenisation process alone is a result of extensive and ongoing research. The algorithm is then trained on this “tokenized” data that may comprise millions or billions of tokenized data elements.

However, copyright protection does not extend to facts, ideas, or mathematical concepts. Although computational analysis may involve technical reproductions of copyright protected works, such reproductions are 1) necessary and incidental to the technical process of training an AI system and 2) do not involve the consumption of copyrighted works for their expressive content. The computational analysis involved in transforming tokenized data derived from on-line content 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. In sum, AI training is not properly regarded as implicating copyright.

It is important to keep these facts in mind when assessing the appropriate responses to the questions raised in the CAIRG Paper, to avoid incorrectly conflating the relationship between training data used as an input to AI systems (regardless of whether the underlying content is copyrighted) and output data which has a direct impact on individuals and may implicate copyright depending on the specific AI generated output.<sup>9</sup>

---

<sup>9</sup> For more information, see US: Artificial Intelligence and Copyright Policy, 08 January 2024 at: <https://www.bsa.org/policy-filings/us-artificial-intelligence-and-copyright-policy>; Hong Kong: BSA Comments on Copyright and AI Consultation, 05 September 2024 at: <https://www.bsa.org/policy-filings/hong-kong-bsa-comments-on-copyright-and-ai-consultation>; Japan: BSA Comments on Draft Approach to AI and Copyright, 12 February 2024 at: <https://www.bsa.org/policy-filings/japan-bsa-comments-on-draft-approach-to-ai-and-copyright>; Japan: BSA Comments Regarding Intellectual Property Rights in the Era of AI, 02 November 2023 at: <https://www.bsa.org/policy-filings/japan-bsa-comments-regarding-intellectual-property-rights-in-the-era-of-ai>; Canada: BSA Submission on Artificial Intelligence and Copyright Policy, 09 January 2024 at: <https://www.bsa.org/policy-filings/canada-bsa-submission-on-artificial-intelligence-and-copyright-policy>; Korea: BSA Comments on AI and Copyright Guidelines, 17 November 2023 at: <https://www.bsa.org/policy-filings/bsa-comments-on-kcc-ai-and-copyright-guidelines>; US: Comments to US Copyright Office regarding Artificial Intelligence and Copyright, 30 October 2024 at: <https://www.bsa.org/policy-filings/us-comments-to-us-copyright-office-regarding-artificial-intelligence-and-copyright>; Singapore: Feedback on Draft Copyright Act 2021, 05 April 2021 at: <https://www.bsa.org/policy-filings/singapore-feedback-on-draft-copyright-act-2021>; Singapore: BSA Comments on Proposed Changes to Singapore’s Copyright Regime, 24 October 2016, at: <https://www.bsa.org/policy-filings/singapore-bsa-comments-on-proposed-changes-to-singapores-copyright-regime>; Australia: BSA Response to Australia DCA Copyright Modernisation Consultation, 07 June 2018 at: <https://www.bsa.org/policy-filings/australia-bsa-response-to-australia-dca-copyright-modernisation-consultation>

## Transparency Regarding AI Inputs

The CAIRG Paper discusses various perspectives regarding transparency and the disclosure of AI training data sets in the context of copyright interests.<sup>10</sup> Specifically, the CAIRG Paper asks whether there are any copyright-related transparency issues that would NOT be addressed by the Proposal Papers' proposed Guardrail 3.<sup>11</sup>

Our response is that **there are no additional copyright-related transparency issues for the CAIRG to consider that are not addressed by the Proposals Paper** with respect to AI input/training data. Instead, **we urge the Australian Government to focus transparency measures in the high-risk AI context and avoid imposing such mandates for copyrighted content.**

As the CAIRG Paper acknowledges, the Proposals Paper includes proposed guardrails related to “transparency” or “disclosure” of AI training data sets. For example, Guardrail 3 proposes the disclosure of data sources to ensure that the data is “obtained legally” and that “[D]ata sets must be disclosed”.<sup>12</sup> In our comments to the Proposals Paper, we agree that training data should be legally obtained but warn against imposing impractical disclosure obligations.

While AI systems vary substantially, large language models and other large AI systems require an enormous set of varied data for training, and it would not be possible to provide a detailed “catalogue” of all the copyrighted works used in the training set in such circumstances. Furthermore, imposing a requirement to disclose data sources could lead to the disclosure of confidential information and trade secrets relating to training methods which involve significant research and development.<sup>13</sup> Measures that undermine trade secret protection would significantly hinder investment in AI development.

Additionally, the Proposals Paper takes a risk-based approach to AI governance and, in general, proposes that the “mandatory guardrails” discussed would be imposed only on developers or deployers in cases where AI systems are used in a way that may result in a “high-risk” to individuals or groups. Whether an AI system is deemed high-risk depends on its intended use, not on the nature of the data used to train the system.

As described above, AI training generally does not harm the legitimate interests of copyright holders and should not implicate copyright protections regardless of whether the publicly available data is copyright protected. Furthermore, the mere use of publicly available data is not determinative of whether an AI system is high-risk. For both of these reasons, the use of publicly available information in AI training does not necessitate additional transparency or disclosure requirements.

To provide legal certainty for both copyright holders and AI developers, **we recommend the Australian Government amend the Copyright Act to adopt an explicit computational data analysis exception**, like those in Japan and Singapore and under consideration by Hong Kong, which can be applied for both commercial and non-commercial activities.<sup>14</sup> Such certainty is necessary to facilitate investments in AI development and deployment in Australia, while clarifying appropriate protections for copyright holders.

---

<sup>10</sup> CAIRG Paper, pages 4-5

<sup>11</sup> CAIRG Paper, page 7

<sup>12</sup> Proposals Paper, Guardrail 3, page 37

<sup>13</sup> BSA Comments, pages 2 and 9

<sup>14</sup> Copyright and Artificial Intelligence Public Consultation Paper, July 2024, at: [https://www.cedb.gov.hk/assets/resources/cedb/consultations-and-publications/Eng\\_Copyright%20and%20AI%20Consultation%20Paper%20\(2024.07.08\).pdf](https://www.cedb.gov.hk/assets/resources/cedb/consultations-and-publications/Eng_Copyright%20and%20AI%20Consultation%20Paper%20(2024.07.08).pdf)

As the Government of Singapore explained when faced with similar questions regarding a potential computational data analysis provision:

A specific exception for such activities is preferred to relying on the general open-ended fair dealing defence, as it promotes certainty and allows calibration of specific safeguards and conditions to address the concerns raised by [various stakeholders]. These include conditions to preserve and protect rights-holders' commercial interests and freedom to conduct business based on licensing and subscription models. For example, the exception will be limited to acts of copying and a user must have lawful access to the works and other subject matter that are copied. If certain material can only be accessed through a paid subscription, the user must pay for the subscription before using the material for text and data mining. The user also cannot distribute the material to anyone without such lawful access. The exception will also not prevent rights-holders from taking reasonable measures to maintain the security and stability of their computer system or network.<sup>15</sup>

In sum, a well-crafted exception would: a) promote AI development and wider economic growth; b) facilitate research; c) allow Australia to maintain its competitiveness on the global stage and sustain its position as a hub for growth and investment; d) afford legal certainty to users; and e) strike an appropriate balance between copyright protection and reasonable use of copyright works.

Beyond the Asia-Pacific region, many other countries have introduced computational data analysis, also known as text and data mining (TDM) exceptions, or have fair use copyright doctrines that support AI innovation. These includes member states of the European Union, Israel, and the United States.

For example, the EU provides legal certainty for AI development by introducing an explicit TDM exception, while at the same time providing rightsholders with the ability to reserve the ability to opt out of TDM activity. In this regard, while AI training does not infringe copyright, and explicit TDM exceptions provide valuable legal certainty, **BSA supports voluntary conversations around automated tools to indicate that the rights-owner does not want a website used for training purposes, similar to the current "do not crawl" tools that apply to search engines.** BSA supports further discussions to arrive at effective, consensus-based technical mechanisms.

The Australian Government has a stated goal of supporting the development and deployment of AI-enabled innovations, a field of critical technology that is in the national interest, and which could contribute \$45 billion to \$115 billion to the Australian economy through the effects of generative AI alone.<sup>16</sup> To achieve these and related objectives, it is important for the Australian Government to promote an intellectual property ecosystem that effectively supports such innovation by providing legal certainty to AI developers and deployers as well as creators that generate creative content that rely on copyright protections. Such an approach will allow Australian industries to seize the many opportunities that AI-enabled innovations present.

## Transparency Regarding AI Outputs

As described in our response to the Proposals Paper, **BSA supports the development and deployment of reliable content authentication and provenance mechanisms (e.g.,**

---

<sup>15</sup> Singapore summarized the limited nature of the exception as follows: (1) The exception will only cover acts of copying (and not other acts protected by copyright); (2) The copying must be for the purpose of data analysis. If no analysis is performed on the work that has been copied, the exception will not apply; (3) Both non-commercial and commercial activities can qualify for the exception; (4) There will be additional safeguards and conditions for the exception to apply. These will include the following: (a) The user must have lawful access to the works that are copied; (b) If access to the works requires payment (e.g., a paid subscription), the user must have paid; (c) The user cannot distribute the works to those without lawful access to the works; and (d) Rights-holders will not be prevented from taking reasonable measures to maintain the security and stability of their computer system or network. Singapore Copyright Review Report (2019), at: [https://www.mlaw.gov.sg/files/news/public-consultations/2021/copyrightbill/Annex\\_A-Copyright\\_Report2019.pdf](https://www.mlaw.gov.sg/files/news/public-consultations/2021/copyrightbill/Annex_A-Copyright_Report2019.pdf)

<sup>16</sup> Proposals Paper, page 5

**watermarking) that can help users identify AI-generated content.**<sup>17</sup> We support efforts by the Content Authenticity Initiative (CAI) to promote the open Coalition for Content Provenance and Authenticity (C2PA) standard for content authenticity and provenance. This standard will help consumers decide what content is trustworthy and promote transparency around the use of AI. In conjunction with watermarking, the CAI approach provides secure, indelible provenance. Embracing open standards like that developed by C2PA facilitates interoperability and enhances the integrity of digital content ecosystems. We acknowledge that what constitutes state of the art in ensuring solutions for content provenance will evolve over time so any governance framework must be designed to accommodate such developments and should provide scope for organisations to assess what is the most relevant solution for them when it comes to content authentication and provenance mechanisms.

Although it appears to require substantial effort on the part of the user to prompt even general large language models to produce copyright infringing content, AI tools can be misused for the commercial dissemination of unauthorized digital replicas of an artist's name, image, likeness, or voice. This improper activity is particularly detrimental to artists who depend upon their reputation and public recognition for their livelihood. As such, while recognising the myriad benefits AI systems are presenting for creators to develop their craft, give expression to their ideas, and develop new artistic works and performances, **we recognize that new tools may be needed to protect artists from the spread of such unauthorised AI-generated digital replicas.**<sup>18</sup>

## Conclusion

We thank the AGD for the opportunity to respond to the CAIRG Paper. We hope that this input will be useful, and we look forward to our continued participation in the CAIRG and other policy discussions regarding the Australian Government's approach to AI governance.

Sincerely,



Jared Ragland, PhD  
Senior Director, Policy – APAC

---

<sup>17</sup> BSA Comments, pages 3 and 9-10

<sup>18</sup> US: Artificial Intelligence and Digital Replicas, 12 June 2024, at: <https://www.bsa.org/policy-filings/us-artificial-intelligence-and-digital-replicas>