

Reply Comments: U.S. Copyright Office Notice of Inquiry and Request for Comments on Artificial Intelligence and Copyright

II. Compelling Disclosure of AI Training Datasets Ignores Practical Realities

There has been significant discussion around compelling AI developers to disclose the datasets on which they have trained. We believe that would be a misguided policy decision for several reasons.

First, such a requirement would be largely unworkable in practice:

- LLMs are generally trained on a wide variety and innumerable quantities of publicly-available source material, such as web-crawled or scraped data, the very nature of which is disparate and ever-shifting. This stands in stark contrast to other technologies that rely on smaller datasets, which lend themselves to being provided or managed by a centralized provider.⁵ Correspondingly, because there is no central set of data vendors or a standardized approach to referring to individual copyrighted work, the copyright status of individual works contained in billions or trillions of datasets are effectively impossible to discern. Furthermore, there is significant variation in the nature and extent of metadata associated with each item of training data, and none that is specific to the copyrighted status of content that is publicly available.
- LLMs are deployed internationally. Consequently, the burden of harmonizing copyright law for any given dataset across geographies prior to deployment would practically mean that no useful model would ever get off the ground.
- Development of foundational models is an ongoing process that involves substantial amounts of work to filter and refine the training data used for each version of a model. Data used for one version of a model may not be used for a subsequent version and vice versa, multiplying the complexity of maintaining an up-to-date list of every copyrighted work that forms part of the training data for a specific version of a foundational model.
- Every iteration of this technology may give rise to new considerations about data use and integration that are impossible to forecast, infinitely complicating the model training process. Consequently, prior to and at various stages during the model training of a new model iteration, the datasets themselves and aspects thereof need to be refined, structured, categorized or manipulated in a multitude of ways to address various factors from model efficacy to safety.

Second, requiring developers to disclose the datasets used to train models would result in the disclosure of proprietary and commercially sensitive information.

Third, and perhaps most important, compelling disclosure of training data for LLMs is at odds with the principle of fair use, a concept that dates to the *Statute of Anne* in 1709 and is a complete defense to copyright infringement under the *Copyright Act*.⁶

⁵ For example, the Bloomberg Terminal is a nearly ubiquitous financial services tool that relies on data of a limited quantum and mainly from selected vended sources.

⁶ 17 U.S.C. § 107. Fair use acknowledges “a privilege in others than the owner of the copyright to use the copyrighted material in a reasonable manner without [their] consent.” *Harper & Row Publishers, Inc. v. Nation Enterprises*, 471 U.S. 539, 549 (1985).

III. Interoperability with Emerging International Standards is Important

Given the global and borderless nature of the development and deployment of LLMs, Cohere believes that it is vital for the U.S. to maintain interoperability with rules, regulations and norms in other jurisdictions for the training and use of AI, such as those recently established or clarified in Japan⁷, Singapore⁸, South Korea⁹, and Israel.¹⁰ In fact, many of these jurisdictions have expressly modeled their laws implicating training of AI to functionally incorporate the concept of commercial fair use under the *Copyright Act*. To the extent that the United States was to distance itself from this long-held principle in copyright law, other jurisdictions would highlight this policy shift to become more attractive hubs of AI training and development. Furthermore, creating an outlier training regime in the U.S. would diminish access to diverse and differentiated AI models, leading to fewer available models in the U.S., depriving American businesses from access to state-of-the-art tools enjoyed by their foreign counterparts. Not only would the number of models be reduced, but the quality of those models would suffer also, as the ability to train on a diverse and global dataset, which is crucial to building AI systems that are safe, accurate, representative, and unbiased across many different regions and cultures, would likewise be significantly diminished.

IV. Technical Measures and Outputs of AI Systems

We are aware of some content creators' concerns that the outputs of AI systems may be similar to the information on which the AI systems were trained. We view such outputs as a solvable technological challenge in AI systems, one which we are committed to mitigating through further development of AI systems and through the implementation of technical measures. We note, as have others, that such phenomena are difficult to reproduce and have been observed mostly through teams seeking to intentionally and actively force a model to generate those specific outputs.¹¹ In other words, other than researchers testing the capabilities, weaknesses, and limitations of early LLMs, users have not reliably experienced reproductions of training data through typical usage. We believe that further testing, evaluation and fine-tuning of foundation models will result in a further reduction of the potential for this to occur. Models are trained to, and fundamentally their value is derived from, their ability to effectively and efficiently provide insights into, and unlock information about, an existing corpus of data, not reproduce it.

V. Addressing Other Points Raised in Comments

We have reviewed a number of the comments provided to the Notice of Inquiry and would like to provide some additional perspective on some of the issues raised in those comments:

⁷ Japan amended its copyright law in 2018 to allow users to analyze in-copyright works for machine learning purposes.

⁸ Sections 190-194 of the Singaporean Copyright Act of 2021 incorporates core concepts of the fair use doctrine from the *Copyright Act* (17 U.S.C. § 107).

⁹ Article 35-5 of the Korean Copyright Act allows for fair use similar to 17 U.S.C. § 107.

¹⁰ Section 19 of the Israeli Copyright Act allows for fair use similar to 17 U.S.C. § 107.

¹¹ See Nicholas Carlini et al., Extracting Training Data from Diffusion Models at 4-7 (Jan. 30, 2023), available at: <<https://arxiv.org/abs/2301.13188>>.

- The suggestion to amend the *Copyright Act*'s Section 1202¹² to remove the “double scienter” requirement would be a mistake and violate international law.¹³ The comment proposes to amend the *Copyright Act* to find liability for those who merely removed content management information, removing the well-established requirement that they must also be aware that such removal would facilitate an infringement. This would be a mistake not only because it would ensnare a slew of innocent and unwitting parties, but it also runs directly counter to U.S. obligations as a party to the WIPO Copyright Treaty and the WIPO Performances and Phonograms Treaty, not to mention the Register of Copyrights.¹⁴
- The suggestion to create new recordkeeping obligations to require AI developers to keep a log of training data, user prompts and other information, along with an ability to access such information through the establishment of a new administrative subpoena process¹⁵, would be a unique and burdensome obligation on the nascent AI industry and put those that employ AI engineers in the U.S. engaged in training AI systems at a strong disadvantage in comparison to their non-U.S. training peers. The proposal also assumes that the AI model is deployed through an application programming interface (API), which in a number of cases for Cohere, it is not.¹⁶ Furthermore, such obligations would also likely not withstand constitutional scrutiny. The net effect of this new requirement would be the creation of more incentives for model developers to train in less burdensome, and potentially less well-regulated, jurisdictions, resulting in models that would be more susceptible to misuse or abuse.
- The suggestion to expand the scope of the current copyright rights under Section 106¹⁷ of the *Copyright Act* to include a specific right to train AI models would be shortsighted and severely impinge long-established and important First Amendment principles.

VI. Conclusion

In summary, Cohere is supportive of an inclusive and informed approach towards the development and regulation of AI. As exciting as this new technology may be, and as optimistic as we are about its future, we do not believe that it is so qualitatively different as to necessitate an entirely new framework to manage it. We believe that the established principles of copyright law – specifically, the fair use (including transformative use) doctrine – are adequate to address issues such as

¹² 17 U.S.C. § 1202(b).

¹³ Comments from Graphic Artists Guild to U.S. Copyright Office available at: <https://www.regulations.gov/comment/COLC-2023-0006-9054>.

¹⁴ According to the Copyright Office, Section 1202's provisions do not and should not apply unless the actor “know[s] or ha[s] reason to know that his acts ‘will induce, enable, facilitate or conceal’ infringement.” WIPO Copyright Treaties Implementation Act and Online Copyright Liability Limitation Act: Hearing Before the H. Subcomm. on Courts and Intellectual Property of the H. Comm. On the Judiciary, 105th Cong. 51 (1997) (statement of Marybeth Peters, Register of Copyrights, Copyright Office of the United States).

¹⁵ See Comments of the American Association of Independent Music and Recording Industry Association of America, Inc. (“RIAA”).

¹⁶ Cohere's LLMs may be deployed through a variety of channels including via an API, multiple cloud vendor marketplaces (available via GCP, Azure, AWS, and OCI), a virtual private cloud or “on premise”. In some of those deployment instances (e.g., if a customer prefers to have Cohere's LLMs deployed through their own data centers), Cohere would have no access to the customer's prompting and related information.

¹⁷ *Id.*



ownership rights, and to appropriately draw the line between permitted and socially beneficial uses of existing works and when permission may be required. We look forward to continuing to work with all stakeholders to clarify the best path forward for all.