**OpenAI**

Position Paper by OpenAI
for IAB Workshop on AI-CONTROL (aicontrolws)
Submitted: August 2, 2024

OpenAI submits this position paper in response to the announcement by the program committee of the IETF's interest in convening the IAB Workshop on AI-CONTROL in September 2024.[1]

OpenAI is dedicated to developing advanced AI technologies to benefit all of humanity. We want our AI models to learn from as many languages, cultures, subjects, and industries as possible so they can benefit as many people as possible. The more diverse datasets are, the more diverse the models' knowledge, understanding, and languages become – like a person who has been exposed to a wide range of cultural perspectives and experiences – and the more people and countries AI can safely serve.

At the same time, we believe AI systems should benefit and respect the choices of creators and content owners. Accordingly, OpenAI pioneered the use of robots.txt as a means to enable web publishers to express their preferences about the use of their content in AI.[2] This approach has since been widely embraced by leading AI developers and webmasters alike. In light of this widespread adoption, it appears that robots.txt has already become, in the words of the workshop invitation, "a de facto solution to AI crawling opt-outs."

Nevertheless, we believe that while robots.txt is an effective mechanism for some content owners, it is also an incomplete solution, as many creators do not control websites where their content may appear, and content is often quoted, reviewed, remixed, reposted and used as inspiration across multiple domains. And we agree with the workshop invitation's observation that "[t]his emerging [robots.txt] practice raises many design and operational questions."

OpenAI has been closely following early investigations regarding various AI opt-out mechanisms. Two recent papers from OpenFuture bear special mention, identifying and grappling with many of the practical challenges that remain to be overcome.[3] In April 2024, OpenAI participated in a workshop on opt-outs convened by the European Commission DG CONNECT Copyright Unit, where some of these same themes were further explored. We also note that the enactment of the AI Act in the EU has created additional interest in this topic, insofar as its Article 53(1)(c) references the text-and-data-mining opt-out provisions of EU copyright law.

---

[1] https://datatracker.ietf.org/group/aicontrolws/about/
[2] https://platform.openai.com/docs/bots
[3] Considerations for Implementing Rightsholder Opt Outs by AI Model Developers, https://openfuture.eu/publication/considerations-for-implementing-rightholder-opt-outs-by-ai-model-developers/; Defining Best Practices for Opting Out of ML Training, https://openfuture.eu/publication/defining-best-practices-for-opting-out-of-ml-training/

One approach beyond robots.txt that has been proposed would rely on embedded metadata standards to signal AI training preferences (TDM-Rep and C2PA represent two examples). However effective these approaches may be with respect to materials created recently enough to take advantage of those metadata standards, they cannot assist with the universe of older materials that make up almost the entirety of publicly available internet materials.

OpenAI believes that more is needed, beyond robots.txt or metadata standards, to fully address the practical needs of AI developers and content owners working to practically and efficiently implement the opt-out preferences of content owners and creators.

As a result, in May 2024, OpenAI announced its intention to develop Media Manager, a tool to enable creators and content owners to identify what they own and specify how they want their works to be included or excluded from machine learning research and AI model training.[4] The goal is to deploy an asset-based content identification system using fingerprinting technologies to identify works to exclude them from AI training. We're collaborating with creators, content owners, and regulators as we develop Media Manager. Our goal is to have the tool in place by 2025, and we hope to rely on interoperable standards that can assist AI developers generally.

We look forward to further exploring these topics with you and learning from others who submit proposals and attend the workshop.

Fred von Lohmann
Associate General Counsel, Copyright
OpenAI

---

[4] https://openai.com/index/approach-to-data-and-ai/