# Feedback on the AI Act

The proposal for a Regulation laying down harmonised rules on Artificial Intelligence (AI Act) is a welcomed effort in the development of specific requirements for "high-risk" systems, prohibited AI practices, and new rules on market monitoring bodies. The AI Act constitutes a crucial step in the pursuit of Europe's strategic autonomy and digital sovereignty, as highlighted in the European Data strategy and AI strategy.

We are submitting these remarks based on ongoing work on the "AI Commons" project, a new research and policy design process aimed at exploring the interconnection between openly licensed content and AI training for facial recognition algorithms. This research and policy puzzle arose following the well-known case of the 2019 IBM "Diversity in Faces" dataset, where more than 100 million openly licensed photos – made available under Creative Commons licenses – were obtained from Flickr and used to train facial recognition algorithms. This AI training dataset has been since then used not only by academia and research organizations, but also by entities involved in the training of AI systems for defense, military and law enforcement purposes. All of this happened without either authors or subjects of the photos being aware that their images were used to train biometric categorization systems. This is due to the fact that open licenses assume that broad permission for use is given, and do not include any consent mechanism.

In the previous consultation on common European data spaces for the Data Governance Act, and on the inception mechanism for the Data Act, we discussed the fundamental role played by Open Access Commons-based data sharing. Our feedback on the AI Act proposal is also based on our interest in defining rules for responsible use of Open Access Commons resources, in particular for AI systems.

Yet, as acknowledged by Creative Commons, this issue of compliance with copyright law is far from being solved since there is currently "no consensus on whether the use of copyright works as input to train an AI system is an exercise of an exclusive right". In other words, the legality of such action strictly depends on the use of said copyright works as AI training material constitutes a reproduction or an adaptation right by rightholders. The latter is nonetheless extremely hard to be determined on an horizontal basis since Creative Commons licenses do not restrict reuse as long as the specific licenses' terms are

respected. Furthermore, Creative Commons acknowledges that setting proper standards or regulation for this case goes beyond copyright matters.

This issue requires greater attention by policymakers as the use of openly licensed content for AI training might collide with ethical and normative considerations which are at the heart of EU data protection law. Besides the extremely broad character of Text and Data Mining exceptions in the CSDM Directive, the GDPR lacks adequate clarity as this issue has not been tested in Courts yet. Recent research suggests that the majority of researchers share concerns over the ethical legitimacy of AI training with open datasets. According to a recent survey commissioned by *Nature*, around two-thirds of respondents indicated that AI training based on facial recognition applications should only be performed after receiving informed consent by users. Yet, as confirmed by the recent "Clearview AI" case, this is hardly the case as more than 3 billion facial images were amassed to train facial recognition softwares, then sold to private companies or law enforcement authorities. On this matter, we strongly support the legal complaints advanced by Privacy International, Hermes Center for Digital Rights, Homo Digitalis and Noyb to the French, Austrian, Italian, Greek and UK data protection authorities.

This issue necessitates further consideration to safeguard citizens' rights to protection of personal data under Article 16 TFEU and Article 8 of the Charter. In this light, the AI Act is crucial to shed light on the relation between users' rights to personal data protection and AI training with openly licensed material. Although the AI Act is not explicitly concerned with the regulation of openly licensed material for AI training, it lays important stepping stones to enhance the protection of users' biometric data.

In the remainder of our submission, we focus on the current definition of biometric data in the AI Act, which might suffer from significant definitional loopholes.

## Explicitly include "faces" in the definition of biometric data

Article 3(33) of the AI act defines 'biometric data' as "*personal data resulting from **specific technical processing** relating to the physical, physiological or behavioural characteristics of a natural person, which allow or confirm the **unique identification** of that natural person, such as facial images or dactyloscopic data*". This definition, which includes facial images in its wording, is nonetheless restricted to data that has already been **technically processed** for use in an algorithmic system. In addition, as the definition applies only when allowing or confirming the **unique identification** of a natural person, emotion recognition and biometric categorisation systems in Title IV could be based on datasets which intentionally and unintentionally avoid the threshold for being considered biometric data. This has a

significant impact on the effectiveness of GDPR safeguards, such as users' consent, as the definition applies only after that data has been technically processed. Likewise, evasive processing practices may be employed to avoid meeting the threshold for data qualifying as biometric. Consequently, the processing and storing of raw biometric data at the foundational stages falls out of the proposed biometric protection standards.

To solve this problem, the European legislator should take inspiration from the definition offered by the "California Consumer Protection Act" (CCPA). Instead of technical processing, this definition targets the ground–level of data collection by focusing on **the ability to extract an identifier template**. Biometric information is indeed defined as "*an individual's physiological, biological or behavioral characteristics* (...) *that can be used, singly or in combination with each other or with other identifying data, to establish individual identity*".

Most importantly, the definition also provides an exhaustive list of biometric identifiers that are used to identify individuals. These are "*imagery of the iris, retina, fingerprint, **face**, hand, palm, vein patterns, and voice recordings, **from which an identifier template**, such as a faceprint, a minutiae template, or a voiceprint, **can be extracted**, and keystroke patterns or rhythms, gait patterns or rhythms, and sleep, health, or exercise data that contain identifying information*". This is an important addition as it significantly narrows down the legislative scope of the proposal to the moment of extraction by shedding light on what would qualify as unique identification. In the AI act, this is instead framed in vague terms and significantly opens dangerous rooms of interpretation on which human features can be regarded as biometric data.

Therefore, the proposed definition in the AI Act, if amended on this line of thought, would significantly enhance users' protection of biometric data as the threshold for legal protection would apply at the early stages of collection while clarifying the scope and meaning of the term "unique identification". Consequently, users' biometric data would be subject to clearer and effective protection, as already stipulated by the GDPR.

August 2021
Open Future Foundation
hello@openfuture.eu