IU Position Paper, September 2023

# Generative AI: Copyright status and recommendations for action

**Authors and performers in fields such as journalism, photography, music, drama, fiction/non-fiction, illustration, visual arts, design etc. are represented by more than 40 associations and unions, which are united under the Initiative Urheberrecht (IU). Europe is home to several million authors of intellectual works, whose creations, services and recordings are the basis for material and immaterial value creation that is indispensable for the European economic community and community of shared values. Since around 2010, their copyrighted works and all activities protected by the related rights, where these are available in digital form, have been used to build large databases for the purpose of training artificial intelligence systems.**
**AI tools are useful in many fields, including the arts, cultural and media sectors, where they are widely applied. However, there is an urgent need to regulate generative AI. The latest thoughts on regulatory solutions, based primarily on expert advice from AI scholars and focused on the AI Act, are presented below.**

As highlighted in the IU statement "Künstliche Intelligenz braucht Leitplanken" ("Artificial Intelligence Needs Guard Rails | Initiative Urheberrecht") of 28 April 2023, there are three levels: INPUT – PROCESSING – OUTPUT. At the INPUT level, a distinction must be made between the selection and acquisition of data (SCRAPING) and TRAINING; this internal differentiation is new to our system of classification. In computer science, PROCESSING (*computation*) is consistently described as a black box; not even the operators of AI systems know exactly what happens during the learning process – *and they do not systematically control it*. The products of generative AI are compiled at the OUTPUT level.

Initiative Urheberrecht (IU) is continuously working on an evaluation of the current copyright status of generative AI systems through exchange between AI researchers from various universities and the Fraunhofer Institute as well as specialist lawyers from associations, unions and collecting societies. This position paper provides an insight into this opinion-forming process, combined with short- and medium-term recommendations for action (see also the supplement entitled "Formulierungsvorschläge für den AI Act" ("Suggestions for the EU Artificial Intelligence Act")).

## TAKING STOCK

**INPUT** consists of two steps: SCRAPING and TRAINING.

Initially, all kinds of data, including significant amounts of copyrighted works and performances that are essential for this first step, are collected and stored so that they can be used to train the AI system in the next step. This process is called **SCRAPING**, and it is undoubtedly a copyright-relevant process. Specifically, the works and performances collected are **stored in a database** so that they can be made available for training.

After training, the data used is no longer directly needed; they can (and should) therefore be deleted. According to leading computer scientists specialising in AI, however, the relevant databases, or rather their contents, are **not** usually **deleted**. There are several reasons for this, not least of which is the possibility of subsequent re-processing to determine comparability.

_____

During **TRAINING**, the second step at the INPUT level, models are taught from the previously stored content that predict probabilities (such as certain character, pixel or word sequences). Current models are generally based on machine learning (including neural networks/deep learning).

In terms of copyright, the processes at the INPUT level can be described as follows:

**SCRAPING** involves mass **copying**, a copyright-relevant process (Section 16 German Copyright Act). All parties involved (rights holders, users, AI platform operators etc.) agree that this is the case from a technical and copyright perspective. All data is **copied**. Scraped data of all kinds is stored in a **database** as a basis for TRAINING.

From a computer science perspective, **TRAINING** itself does NOT result in a usable database "per se"; the computed *model* cannot and is not intended to function like a traditional database. The PARAMETRIZATION of the trained data results in a highly abstract REPRESENTATION or MANIFESTATION of the content within the model. From a copyright perspective, there is no database in the sense of Art. 1 of the Database Directive, since the data is not "individually accessible by electronic or other means".

After training, the data in the model is not available as copies in the "traditional" sense. The AI model no longer uses the database created previously for its output, but rather only uses the parameters selected from that database. Based on the current state of the art, however, it is not possible to provide an unambiguous description of how exactly the parameters in the model are classified, even from a technical point of view. In the copyright debate, the question of whether reproductions in the sense of copyright law still exist after the training has been completed is the subject of some debate. However, there is much to suggest that even the trained AI model (at the 2nd level) still contains **reproductions** in a copyright-relevant sense, since it is undoubtedly possible for systems like ChatGPT to reproduce poems or other copyrighted texts. Even if the reproduction of the respective text is based on the probability of stringing together the respective passages based on the respective user requests ("prompts"), the work is still part of the model in this way. Both German and European copyright law define reproduction broadly and **independently of technology** (see Art. 2 InfoSoc Directive: "…direct or indirect, temporary or permanent reproduction by any means and in any form…"). This includes the transmission of parameters that allow the work to be reproduced (albeit by means of probability calculations or the like).

**It is by no means clear that the reproductions made during scraping for use in massive machine learning algorithms and the creation of foundation models are covered by the legal permission for text and data mining under Section 44b of the German Copyright Act;** moreover, the purpose of TDM is not to generate new content, but to explore the data. While it is true that data analysis takes place in the training of generative AI systems, "patterns, trends and correlations [...]" is not *obtained* "for the purpose of gathering information" as described in Section 44b of the German Copyright Act, but the features obtained are "internalized" – they are neither understandable nor accessible to humans. The form of the content, rather than the content itself, is thus represented in an abstract way; one can imagine this as the categorical difference between a package insert and a drug or a recipe and a dish. Accordingly, when training generative AI systems, the knowledge gain of the TDM provision is not the primary focus.

We assume that the **description of TDM in Section 44b does not correspond to what actually happens at the INPUT level in the collection and processing of works for machine training/learning**, but acknowledge that there are contrary opinions. However, should the view prevail that the processes described above are covered by the TDM provision, **a remuneration obligation for the uses occurring is imperative and must be established immediately.**

_____

The assumption that the legislator did not intend for the current TDM exception to allow AI scraping and training is all the more plausible given that MEP Axel Voss reported at the Erich Pommer Institute Copyright Conference in June 2023 that AI had not even been considered when the TDM exception was introduced. **Unlicensed copying and other copyright infringements occur constantly at the INPUT level.**

As noted above, a significant amount of the scraping that has occurred to date took place well BEFORE the TDM exception in the DSM Directive came into effect in 2019. Accordingly, the appropriation of these vast datasets without consent, attribution or remuneration cannot be legitimized by reference to TDM under any circumstances. These are blatant copyright violations that Germany and the EU cannot accept, if only for economic reasons; after all, the content appropriated has been used to train systems that are preparing to replace the commercial production of new works and recordings by authors.
It is essential to find solutions for past use of this data that are acceptable to the rights holders.

**Unlearning**/forgetting what has been learned is not possible according to the current state of technology and statements by leading AI scientists. There is therefore a risk of substantial claims for damages. In the US, there is talk indicating that if one of the pending lawsuits against generative AI providers is successful, their entire MODEL would have to be deleted and the training process would have to be restarted.

**Proof** of whether certain specific works have been used for machine learning or for the creation of the foundation models cannot be provided at the INPUT level alone, but must often be provided at the OUTPUT level, in addition to the transparency we require about the type and amount of training data used: **For example, if a prompt asks for the style of a particular artist, and the output is very close to that style** ("proximity"), then it can be concluded that the works of that artist were used for training.

If the work is no longer present as a copy in the *model*, the corpus of training results, but is represented abstractly, then this may constitute a **NEW TYPE OF USE**. If it can be proven at all that the specific work exists and can be found, the question of whether it manifests itself in abstract vectors or in bits and bytes is irrelevant. We are dealing with a **technology** that **allows REPRODUCTION**.

In many cases, no pixel in the OUTPUT product, such as images, is identical to the original, which raises the question as to whether proximity should be determined technically or based on reception.

**OUTPUT** level

The assessment of the OUTPUT, i.e. the products of generative AI, is essentially unchanged. Copyright protection requires "individual intellectual creation" associated with a natural person. As this is not the case with autonomously generated AI products, these types of products cannot be given the status of a work and therefore cannot be given copyright protection. Nor can the person formulating the prompts claim any rights with respect to the result on the basis of the prompts alone, because the mere formulation of the task and the choice between several results proposed by the AI system is not a creative act.

_____

The situation is different if the AI is used merely as a tool, possibly even as one of several tools, in which case the creative act in question is likely to lie with the author making specific use of the technology, provided that they are a natural person. Collecting societies, stock photo agencies and other organizations that manage large repertoires will need to develop strategies for dealing with such works and performances and adapt their rules accordingly. If the output is a work that (still) falls within the scope of protection of a pre-existing work, it will be even less possible than usual to speak of the kind of "parallel creation" that is theoretically possible under copyright law. A permissible parallel creation presupposes that the author was not aware of the older work. This cannot be assumed in the case of AI models due to the data used to train them.

The following aspects are also relevant:

Generative AI system output is based on the training that has taken place – and thus on the content used in the course of the training. Therefore, the output generated cannot be viewed and evaluated in isolation from the input. If the content used in the course of the training is clearly recognizable in the output, comprehensive regulations are required for its handling, which in our view can be reduced to the formula "**3C+1T**": *Consent/Credit/Compensation + Transparency*, where the first three are impossible without transparency on the part of the AI provider.

This formula already addresses the **need to protect the personal rights** of all potentially involved and affected parties. When voices are separated from people in voice cloning, when actors and actresses are replaced by their own clones, when inputting an artist's name into a text-to-image model returns countless products that give the impression that they were created by that very artist, but also when the protagonists of journalistic and documentary media have statements put into their mouths that they never made and would never make, this is a profound encroachment on the personal rights of those affected. The current actors'/screenwriters' strike in the USA is proof of the urgency of this aspect. There is a need for clear and enforceable rules, including a right of prohibition, to protect the personality, to which – especially in copyright law – the livelihood of most stakeholders is closely tied.

Any attempt to use contractual agreements to allow the unrestricted use of input for the production and operation of systems that compete directly with the authors of the training content should be prohibited entirely.

No agreement has been reached within the IU on possible ancillary copyright claims on the products of generative AI systems. However, it should be borne in mind that intellectual property rights, such as those of the sound carrier or film producer, are based on the idea of investment protection. Providers of AI services use the AI infrastructures of companies such as Microsoft, for example, meaning that the investments that may be worthy of protection are essentially not made by the individual providers.

From the point of view of the IU, it is essential to leave the legal link between author and copyright untouched and not to grant copyright protection to autonomously produced AI products.

For the sake of differentiation, and in view of the danger of manipulation and misinformation, AI products should be labelled clearly and comprehensively – automatically from the moment of their creation, if necessary. For example, the ISCC standard[1] developed with EU funding could be helpful,

---

[1] Information about the ISCC can be found at: https://iscc.codes /// ISCC currently has the status of "Draft International Standard" ISO/DIS 24138; ISO project page: https://www.iso.org/standard/77899.html

especially since it is decentralized and non-proprietary. The deletion of such a label and the separation of any metadata record potentially associated with the file or its contents from the file or its contents should be prohibited – similarly to the prohibition of circumvention of copy protection measures.


## RECOMMENDATIONS FOR ACTION

Urgent action is required, given the rapid development and proliferation of the technology in question. We therefore expect the EU's AI Act to do the following:

- Introduce a **comprehensive transparency obligation** that, in addition to its direct copyright relevance, also allows for market monitoring and clear impact assessment. Authors, performing artists and rights holders must be able to find out whether and to what extent their works and performances are being used for training at the INPUT level and the extent to which they are being used as a basis at the OUTPUT level.
- Introduction of a fundamental duty to **label products originating from generative AI,** facilitating the unambiguous, comprehensible identification of machine-generated content. However, it is possible that a total and comprehensive duty to label may not apply in certain rare cases due to constitutional requirements.
  - Incidentally, labelling would also be in the interests of AI providers who need to strictly avoid feeding their systems with AI-generated products as training data in order to prevent a *model collapse* or the development of a *degenerative AI system*.
- Human rights must be reserved for humans, as must copyright protection. As demanded by the European Parliament, **fundamental rights** and **copyrights** must be respected.
- Copyright protection, as a property right, is based on fundamental rights. In addition, the aforementioned copyright infringements have fundamental rights **implications for the personal rights** outlined above.
- **Proof** and **liability** must be clarified in the AI Act and, if necessary, the present text of the Act must be supplemented.
- **Contractual provisions** on the unrestricted use of performances via AI **must be prohibited**.

Since we are convinced that the courts will confirm our view that there is no legal basis for scraping and training at least, we would like to point out once again that this letter is not intended to withdraw any of our demands, but rather to limit our requests in a pragmatic and solution-oriented manner to demands that are specifically directed at the AI Act.

Therefore, for the time being, we are only talking about inclusion in a regulation that is intended to continue to apply. We will work hard for the following demands:

- Should the acquisition and use of works and recordings not be covered by the TDM exception, the instrument for a legally secure way to avoid long and unpleasant disputes in court is **licensing**. All stakeholders on the rights holders' side of the market are prepared to enter into solution-oriented licensing negotiations.
- All uses, including those prior to entry into force of the DSM Directive, must **provide for adequate remuneration**. In view of the perpetual value of the use made, the remuneration must be **substantial,** and the remuneration obligation must be of a **long-term** nature.
- Essentially, the legal framework for TDM needs to be clarified, corrected, concretized and made subject to compensation. There are several possible models.

---

- Regardless of the preferred remuneration solution, the remuneration of the original rights holders (authors and artists) must be guaranteed. Remuneration for rights holders is not necessarily synonymous with remuneration for authors.
- An opt-in, in the spirit of copyright law, is preferable to an opt-out, which is not in the spirit of copyright law per se in order to preserve the decision-making option, also in relation to the author's personal rights – and thus the right to say "no" to harmful conditions and uses; this is only possible in the absence of an exception.
- A statutory license with an opt-out option can both enable authors to object individually to the use of their works and at the same time ensure that authors who have not opted out are adequately remunerated.
- If, however, the TDM exception is applicable, remuneration must be made mandatory and a practicable opt-out option must be provided for authors.
- Contrary to the currently applicable Section 44b (3) sentence 2 German Copyright Act, machine-readability may only be required if standards are already in place based on which authors can formulate their reservations, and if authors have the possibility to sanction violations.

In its new data strategy, "Fortschritt durch Datennutzung" ("Progress through Data Utilization") published a few days ago, the German government also reiterated the need to protect intellectual property and other fundamental and property rights.[2]

Overall, the requirement adopted by the German Ethics Council should apply: in legal terms, **HUMAN CREATIVITY and achievement** should be valued differently and more highly than their machine imitations. The German Cultural Council also follows this thinking in its statement.[3] Politicians must consider that the value creation of the entire national and European creative industries takes place and is accounted for locally, while the profits generated by AI providers – together with the cultural heritage, world knowledge, innovative power and identity-forming personal intellectual creations of all European knowledge workers in their entirety – are not realized in the EU, but in the USA and China.

Science journalist and author Ranga Yogeshwar describes the current situation as follows:
"[…] *we are currently experiencing the greatest theft in human history. The richest companies in the world, such as Microsoft, Apple, Google, Meta or Amazon, are seizing the sum total of human knowledge. That is, all texts, artworks, photographs etc. that exist in digitally exploitable form, in order to then wall off this world knowledge in proprietary products. There is no clear disclosure of what learning data they are using to train the AI in the process. […] Copyright is being ignored – deliberately. Meanwhile, plagiarized products can be produced en masse via AI, with entire professions facing their existential end.*"[4]

**We are currently preparing a statement including legal elaborations and proposals is in preparation; we refer again to our supplement "Formulierungsvorschläge für den AI Act" ("Suggestions for the EU Artificial Intelligence Act").**

**Berlin, 19 September 2023**

---

[2] https://www.bmwk.de/Redaktion/DE/Publikationen/Digitale-Welt/fortschritt-durch-datennutzung.html

[3] https://www.kulturrat.de/positionen/kuenstliche-intelligenz-und-urheberrecht/

[4] Augsburger Allgemeine, 17 May 2023, https://www.augsburger-allgemeine.de/wirtschaft/ranga-yogeshwar-interview-ueber-ki-der-groesste-diebstahl-in-der-menschheitsgeschichte-id66385936.html

*Initiative Urheberrecht represents the interests of approximately 140,000 authors and performing artists in the fields of fiction and non-fiction, visual arts, design, documentary film, film and television, photography, illustration, journalism, composition, orchestra, drama, game development, dance and many more.*

Initiative Urheberrecht

Markgrafendamm 24, Haus 18  |  10245 Berlin

+49 (0)160 9095 4016  | www.urheber.info

Katharina Uppenbrink, Managing Director, Initiative Urheberrecht
katharina.uppenbrink@urheber.info
Matthias Hornschuh, composer and spokesman for IU creatives in the Initiative Urheberrecht
matthias.hornschuh@urheber.info