

# Kluwer Copyright Blog

## Are AI models' weights protected databases?

Nuno Sousa e Silva (Universidade Católica Portuguesa) · Thursday, January 18th, 2024

The ongoing Artificial Intelligence (AI) revolution has machine learning models at its core. Contrary to classic computer programs written by developers, many of these models rely on vast artificial neural networks trained in giant amounts of data. In general, they use what is called a [transformer architecture](#). No one individually writes or encodes these models; they are generated through an automated process of training. Once the data has been prepared and the architecture has been defined (the features and characteristics of a model defined prior to training are called hyperparameters), the computers will run for a long time and with a high cost in order to acquire knowledge “on their own”. The final result, i.e. the model, consists of two files – a simple run file that establishes the functioning of the model (the model architecture) and a much larger file of parameters or weights (expressed as [floating point numbers](#)). Weights are a mathematical expression of the connection between the neurons that make up the network. As [Martin Andreson puts it](#) “*In machine learning research, the weights are everything – the ultimate ‘gold’ that emerges after weeks or even months of training a system.*”



Photo by Steve Johnson on Unsplash

Much has been written about the legal challenges and qualifications of the training process (knowing whether it's legally permissible to train these models on copyright-protected material) and the outputs of these models (especially if there is copyright in the results generated) see [here](#), [here](#), [here](#) and [here](#). However, it's unclear whether *the models themselves* are currently protected by intellectual property laws. The run file is a classic piece of software and does not pose particular difficulties. It's the parameters or weights (the larger file) that raise the puzzling questions.

Including model weights in the realm of copyright faces some obstacles. First, the weights are

numerical values, not exactly an expression of software through code. In other words, there is not a textual production in a programming language.

Second, no programmer actually controls the generation of the expression. Of course, it's possible to argue that there is not an inherent limitation on the kinds of expressions of computer programs (after all, a text in an invented language [but **not the language itself**] can still be protected by copyright) and that compiled code is numerical (binary code), rather than textual, also unreadable by humans. On the topic of authorship, it could be argued that someone (one person or a group of persons) will still have high-level control of the training process, and that control could be enough to attribute authorship. Naturally, it's debatable whether that creates enough of a causal link for authorship (Ginsburg & Budiardjo would probably consider them to be "authorless") – very often the outcome is not only unpredictable but also **uninterpretable**.

A third challenge in granting copyright in the machine learning model weights is their functional nature. This one seems harder to overcome. In fact, each weight is a simple instruction, which can only be expressed that way, i.e. as a numerical value. Hence, following the merger doctrine stated in C-393/09, BSA, §49: "*where the expression of those components is dictated by their technical function, the criterion of originality is not met*". For these reasons, authors such as Hao-Yun Chen, Peter Slowinski, and Begoña Gonzalez Otero seem to reject the protection of models under copyright law.

However, in the EU there is another strong candidate for protecting model weights: the *sui generis* protection for databases established in Directive 96/9. The *sui generis* right is granted (only) to EU-based companies and individuals (Article 11) that make a substantial investment in either the obtaining, verification or presentation of the contents of the database. According to recital 17 of the Directive: "*the term 'database' should be understood to include literary, artistic, musical or other collections of works or collections of other material such as texts, sound, images, numbers, facts, and data; (...) it should cover collections of independent works, data or other materials which are systematically or methodically arranged and can be individually accessed*" (our emphasis). Recital 23 also clarifies that "*the term 'database' should not be taken to extend to computer programs used in the making or operation of a database*".

When looking to apply the EU *sui generis* database protection to machine learning models' weights, one needs to establish: 1) that these can qualify as a database and 2) that there is substantial investment in them by the database maker.

### **Can model weights qualify as a database?**

Earlier case law, C-444/02, OPAP seemed to indicate that the notion of a database requires an index function. As the Court put it "*classification as a database is dependent, first of all, on the existence of a collection of 'independent' materials, that is to say, materials which are separable from one another without their informative, literary, artistic, musical or other value being affected.*" (§29) and "*Classification of a collection as a database then requires that the independent materials making up that collection be systematically or methodically arranged and individually accessible in one way or another. While it is not necessary for the systematic or methodical arrangement to be physically apparent(...) the collection should be contained in a fixed base, of some sort, and include technical means such as electronic, electromagnetic or electro-optical*

*processes (...) or other means, such as an index, a table of contents, or a particular plan or method of classification, to allow the retrieval of any independent material contained within it* (§30). **Considering this understanding, it may be difficult to ascertain if there is a “systematic or methodical arrangement” in model weights.**

On the other hand, one can easily find values (weights) in a parameter file using a simple search function. Nevertheless, those numbers will be meaningless from a human perspective. **More recent case law has adopted a broader notion of a database.** In *C-490/14, Verlag Esterbauer*, it was held that a topographical map can qualify as a protected database. The court stressed “the intention of the EU legislature to give broad scope to the definition of the term ‘database’” (§26). The total model weights will still be a collection of independent numerical values. Furthermore, “the autonomous informative value of material which has been extracted from a collection must be assessed in the light of the value of the information not for a typical user of the collection concerned, but for each third party interested by the extracted material” (§27). This seems to mean that if there are people who see value in the materials (namely the model weights) there will be a database. **In that light, I believe there is enough room to qualify the content of a model weights file as a database.**

### **Is there a substantial investment in model weights by the (database) maker?**

On the requirement of investment, the Court of Justice established that “*The expression ‘investment in ... the obtaining ... of the contents’ of a database in Article 7(1) of the directive must be understood to refer to the resources used to seek out existing independent materials and collect them in the database. It does not cover the resources used for the creation of materials which make up the contents of a database.*” (*C-203/02, BHB v William Hill*, §42 and *C-46/02 Fixtures*, §49).

**This requirement is hard to satisfy when it comes to the model weights.** The amount of money and effort invested in the training of the models is beyond dispute; it’s **the nature of the weights that is a challenge.** While it’s true that those elements do not pre-exist – the weights are the output of the model training – they are a representation of pre-existing information that is **somehow stored in that format.** Another way of looking at it is saying that as the weights are nothing but numbers, they pre-exist (not only as abstract numbers, but also specifically in the data that is being used to train the model), and the investment in the model training is directed to collecting and organizing those pre-existing numbers, not at creating them. *Matthias Leistner* writes that “*the investment intensive methodical or systematical structuring of raw data might be covered under the head of investments in the presentation of the contents of the database.*” **I submit that the definition of the weights that take place in model training might very well qualify an intensive, methodical, and systematical structuring of the dataset being used for that training.**

### **If model weights qualify as a database, what then?**

If we assume that the database rights will apply to model weights, this protection will be limited to EU-based persons and will be granted against extraction and reutilization of a substantial part of the database (for a recent analysis of CJEU case law on the scope of protection, see [here](#)). What this means in practice is giving model developers another tool for controlling the distribution and

use of the parameters file.

The *sui generis* right has a term of 15 years. but it can be renewed with any “*substantial change, evaluated qualitatively or quantitatively, to the contents of a database, including any substantial change resulting from the accumulation of successive additions, deletions or alterations, which would result in the database being considered to be a substantial new investment, evaluated qualitatively or quantitatively*” (Article 10(3) Directive 96/9/EC). Therefore, a model retraining or finetuning, changing the weights, could lead to a new term of protection. The tricky question would be whether that’s still the same database or a new one.

If and when the retraining or finetuning is done by a person other than the rightholder, that could qualify as a derivative database. As long a substantial part of the original database is still used in the derivative database that will count as re-utilization under Article 7(1((b) of the Directive. In other words, the *sui generis* right would allow control of downstream uses of a model.

While model piracy (in the sense of misappropriation of model weights)[1] may not be of high concern at the moment, this protection would provide a higher degree of control to the developers of models. Furthermore, it would give some teeth to open-source licensing conditions, providing means of enforcement beyond contract law. However, as noted, it will only apply to EU-based companies and individuals, which could be perceived as an unwelcome divergence from international standards.

---

[1] Note that the [weights can be shared](#) and implemented using different run files if they have the same architecture. This means that the model can be transferred without using the original copyright-protected run file. [Transfer learning](#) is also possible, i.e., using a model trained for one task as a starting point for a different but related task.

---

*To make sure you do not miss out on regular updates from the Kluwer Copyright Blog, please [subscribe here](#).*

## Kluwer IP Law

The **2022 Future Ready Lawyer survey** showed that 79% of lawyers think that the importance of legal technology will increase for next year. With Kluwer IP Law you can navigate the increasingly global practice of IP law with specialized, local and cross-border information and tools from every preferred location. Are you, as an IP professional, ready for the future?

Learn how **Kluwer IP Law** can support you.

---

79% of the lawyers think that the importance of legal technology will increase for next year.

**Drive change with Kluwer IP Law.**

The master resource for Intellectual Property rights and registration.



2022 SURVEY REPORT  
The Wolters Kluwer Future Ready Lawyer  
Leading change

This entry was posted on Thursday, January 18th, 2024 at 8:03 am and is filed under [Artificial Intelligence \(AI\)](#), [Authorship](#), [Copyright](#), [Originality](#)

You can follow any responses to this entry through the [Comments \(RSS\)](#) feed. You can leave a response, or [trackback](#) from your own site.