

Protecting Artificial Intelligence Innovations as Intellectual Property: Opportunities and Pitfalls

Charles R. Macedo, Amster, Rothstein & Ebenstein LLP, with Practical Law Intellectual Property & Technology

Status: **Maintained** | Jurisdiction: **United States**

This document is published by Practical Law and can be found at: us.practicallaw.tr.com/w-026-6317

Request a free trial and demonstration at: us.practicallaw.tr.com/about/freetrial

An article discussing challenges and opportunities for protecting key components of artificial intelligence, particularly machine learning technology, as intellectual property, including through patents, copyrights, contracts, trade secrets, and other technological mechanisms.

While the basic technological tools associated with Artificial Intelligence (AI), including Machine Learning (ML), date back decades, the commercial application of this technology has exploded with the advent of more powerful computers and large datasets. Research and practical applications of ML in particular continue to be developed and commercialized every day. However, the legal frameworks traditionally used to protect computer-implemented innovations, including patents and copyrights, make it difficult to protect and defend ML innovations.

This article:

- Provides a high-level explanation of key components of ML technology that organizations and their counsel should consider for potential protection and monetization.
- Analyzes the challenges and opportunities of protection available to ML technology through patents, copyrights, contracts, and trade secrets.

This article assumes that, at least for now, patent and copyright protection will not allow computers to be inventors or authors and provides alternative strategies to protect and monetize ML innovations and their valuable components.

For a discussion of key issues concerning IP protection of Artificial Intelligence generally, see [Practice Note, Artificial Intelligence Key Legal Issues: Overview: AI as Intellectual Property](#). For a collection of cross-practice legal issues concerning AI, see [Artificial Intelligence Toolkit](#).

ML Components to Consider for IP Protection and Monetization

AI generally refers to the theory and development of computer systems that perform tasks that traditionally require human intelligence, such as visual perception, speech recognition, decision-making, and translation between languages. ML is a widely-used form of AI today in which a computer system learns automatically from experience. Typically, ML systems build algorithms based on machine learning models that are applied to large data sets.

At a high level, an ML system can be thought of as one that uses a “black box” to provide a desired output in response to an input data query. More specifically, the ML software used to implement an ML system is a black box because, unlike a traditional heuristic (or human programmed) software code, the actual logic (or ML algorithm) the software applies is generally unknown.

Instead, computer engineers know how the software (and the ML algorithm) is **built**, including:

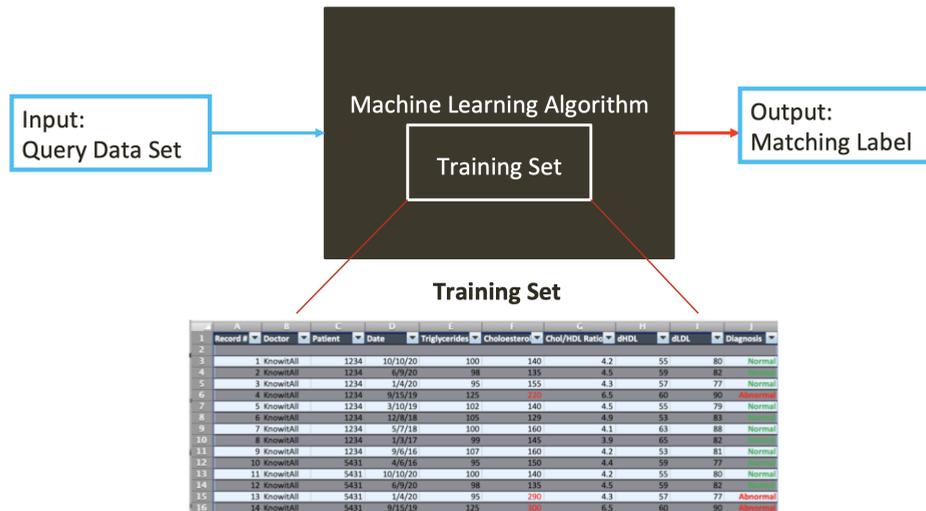
- The training set of data samples, including the features, their labels and “tags” (see Training Set).
- The specified form of an input query data set (see Input Query and Output).
- The specified form of output matching an identified label based on the tags (see Input Query and Output).
- The model upon which the machine learning algorithm is based (see ML Algorithm).



Protecting Artificial Intelligence Innovations as Intellectual Property: Opportunities and Pitfalls

These four components inform how the black box algorithm is “built” and can be used to define a particular ML system, even if the actual logic implemented is never known, as shown here:

Together, these components may be used in a larger system, which can provide more opportunities for IP protection (see Larger System Incorporating Machine Learning).



Training Set

The training set is composed of a large number of data or image samples that include:

- Common **features** (shown in the columns labeled “Data 1,” “Data 2,” above). The features are the data items the ML system uses to identify and compare the data.
- **Tags** that usually reflect a qualitative judgment as to the meaning of the sample data set, such as “Normal” and “Abnormal.” The tags are the result of the analysis of the features that the system uses to generate a useful output.

The key concept is that a training set includes many samples, and each sample includes features with labels (typically the items that are being tested) and one or more tags (typically reflecting the results of what was traditionally a human analysis).

For example, in a training set involving medical data, each item may reflect:

- A set of test results (features, such as blood tests) taken of a specific patient at a specific time.
- The doctor’s tag as to the meaning of the blood test result (for example, normal, abnormal or perhaps a diagnosis like diabetes or heart disease).

Hypothetical Blood Test Training Set

A hypothetical ML system concerning medical data could be used to analyze a blood sample taken from patient on an identified date to test A1C blood levels (a measure of blood sugar concentration and indicator of diabetes), for example. A result exceeding an identified threshold, “210” in the figure above, would be outside the recommended parameters, resulting in a diabetes diagnosis. The patient’s sample would be reflected in one row of the table that could potentially include thousands or millions of rows reflecting other test samples. Each row could include:

- The “Patient ID,” like a date and patient number (to keep the patient’s data anonymous).
- The “A1C” test results, which here would be “210.”
- The resultant diagnosis of “Diabetes.”

In the figure above, the table includes a header with column labels that identify the reflected features and tags. Here, the features are the test results for A1C. The tags are the diagnosis, such as “diabetes.”

This simplified hypothetical includes only one feature, tag, and patient sample. In practice, there could be many features making up what is often called a “feature vector,” and can be one or more tags in a training set, depending upon the qualities analyzed. Further, to be

useful, a training set will need to include a large number of samples, ideally thousands or perhaps even millions of samples, in order to provide enough information to make the training set statistically useful. Also, any samples included that lack a useful tag will be worthless as just “noise” because the sample will not be determine whether the data reflects, in this example, whether the patient has diabetes.

Training Set Characteristics

The training set in the above hypothetical is based on data (the numerical test results for A1C blood level). However, training sets may also be based on images, videos, or sound files, instead of alpha-numeric data. In the simplest case, each image, video, or sound file is a sample that includes at least one tag.

For example, a sample may be an x-ray taken of a patient on a given date, and the tag may be the diagnosis of “cancer” or “no cancer.” The training set would then include just the image and the tag, or may also include data associated with the image, such as the patient’s information, test date, and perhaps other test results.

Training sets can also be more complex and may be made up of a combination of data items or images, videos or sounds, or the like for each sample.

Input Query and Output

When defining a ML system, it is important to understand both the form of:

- The query that will be input into the black box algorithm (and tested against the training set).
- The output resulting from the black box operation.

The query is a specific collection of information that is being tested by the ML system. It will typically consist of one or more features included in the training set (for example, the blood tests run against the medical data training set discussed above). If the features of the query do not match the features in the training set, the designer or developer will need to normalize the query or no useful result will be available.

The output in turn will typically match the tags of the training set. For example, if the blood tests are tagged

with a diagnoses of “normal” or “abnormal,” then the output of the ML system will typically also be “normal” or “abnormal.”

ML Algorithm

The ML system will need to use some form of algorithm in order to allow the input query to be tested against the training set and identify an output. The algorithm is generated using an ML model, provided by a programmer, which the system applies against the training set. When the ML model is applied against a training set, the computer will automatically generate the ML algorithm which is generally considered an unknown “logic” to test the query. So, while a computer engineer will probably not know what the ML system’s applied logic will be, it is possible to know:

- The training set used (including its features, labels, tags and samples).
- The ML model applied against that training set.
- The form of query used as an input and the form of output specified.

As discussed below, each of these components may by themselves be valuable and worthy of some form of IP protection and monetization.

Larger System Incorporating a Machine Learning System

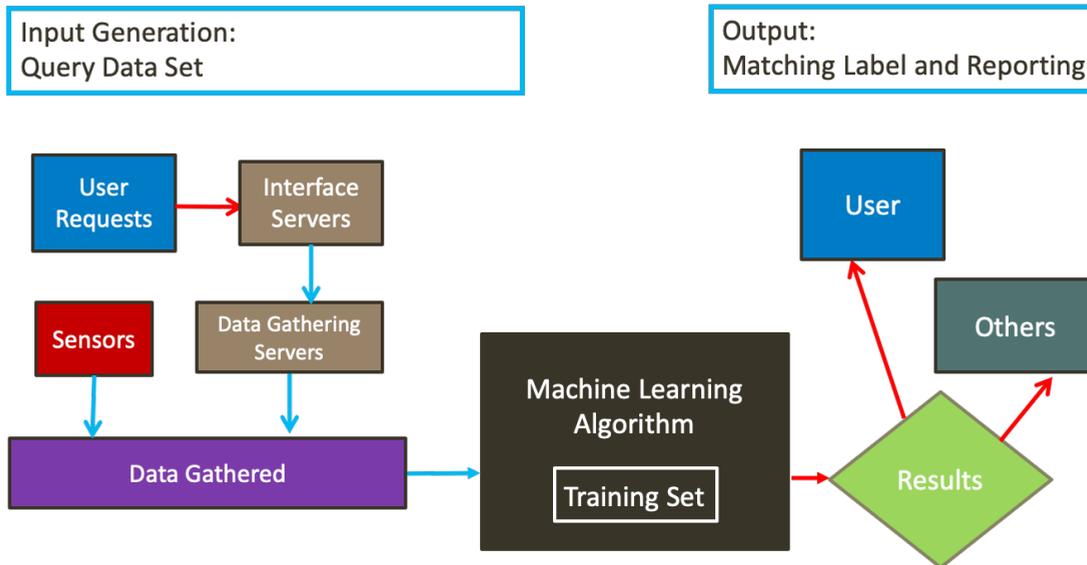
A ML system does not stand alone, but is usually part of a large system that involves many additional processes and systems. In addition to each of the building blocks of the ML system, innovation and value may be found in how such a system fits into a larger system, such as the blood sample testing system described above.

For example, innovation may lie in:

- How the query is generated.
- The processes used to gather the data elements or images to be tested.
- What the system does with the output of a query after it is generated.

Protecting Artificial Intelligence Innovations as Intellectual Property: Opportunities and Pitfalls

An example of such larger system is shown here:



Input Generation as Part of a Larger System

The larger ML system may use a robust series of systems and processes to generate the input data set that the ML algorithm will test. For example, the system may require a user to submit a request through a mobile application or website hosted by interface servers. The interface servers will provide data gathered by the request, which the system may combine with data obtained elsewhere, such as with sensors, to generate the data gathered and in turn generate the query data set.

Output as Part of a Larger System

This larger system also may include series of systems and processes that act upon the output of the ML algorithm. In a simple case the system may simply report the output to the user and perhaps others. In a more complex system, the larger system could act on the output to take action based on the results or trigger a series of other processes that depend on the result obtained. For example, after analyzing an X-ray through a ML process, the system could send the resulting diagnosis to a treating physician with the X-ray to determine next steps.

Challenges and Opportunities for Seeking ML IP Protection

When determining where to protect and monetize ML innovation, innovators should identify where the value best lies in the technology. These value points are the

points where the provider is typically providing something unique or valuable and worthy of protection and compensation. Organizations and their counsel should consider the potential value points in an ML system discussed above, including:

- The training set (see Protecting Training Sets).
- The ML algorithm or model (see Protecting the ML Algorithms or Models).
- The input query (see Protecting the Input Query and Processes Used for Generation).
- The output result (see Protecting the Output Result).
- The larger system incorporating the ML system (see Protecting the Larger ML System).

Patents, copyrights, contracts, and trade secrets are potentially available to protect and monetize computer implemented ML innovation. Unfortunately, the laws and legal paradigms associated with each of these forms of protection, which were developed in a brick and mortar world, present challenges when a computer or non-human inventor or author is involved.

For example, courts and the US Copyright Office agree that a monkey cannot be an author under the Copyright Act (and so presumably neither could a computer) (see *Naruto v. Slater*, 888 F.3d 418 (9th. Cir. 2018)). The US Patent and Trademark Office (USPTO) has also rejected at least one patent application naming a computer system,

DABUS, as an inventor on the grounds that only humans can be inventors (U.S. Patent and Trademark Office, Decision on Petition In re Application of Application No. 16/524,350, available at https://www.uspto.gov/sites/default/files/documents/16524350_22apr2020.pdf).

Until the law catches up with technology and development, it is important to focus any plans to protect and monetize AI/ML innovation on the human-created elements. For more on AI IP ownership considerations, see [Practice Note, Artificial Intelligence Key Legal Issues: Overview: AI IP Ownership](#).

Protecting Training Sets

With the development of very large data sets in many industries, the potential to monetize data as a training set is larger today than ever. Protection may be available through copyright, patent, and trade secret, as well as contracts.

Copyright

Copyright law may in some cases be available to protect training sets, but not always. This is because copyright law protects the original “expressions” of ideas, not ideas themselves. Therefore, while some datasets may be protected under copyright law when they include sufficient original human expression, other data sets may be unprotectable if they lack sufficient human expression (compare *Feist Publications, Inc. v. Rural Tel. Serv. Co.*, 498 U.S. 808 (1990) with *Key Publications, Inc. v. Chinatown Today Publishing Enterprises Inc.*, 945 F.2d 509 (2d Cir. 1991)). For example, a training set may be protectible where expressive thought is used in selection process of the data samples, the categorization and organization of the dataset, or the tagging of the samples.

In order to increase the likelihood of successfully registering and protecting a database under copyright law, the US Copyright Office recommends altering the structure or content of the training set to incorporate greater creativity (see *Compendium of U.S. Copyright Office Practices*, Third Edition, released on 22 December 2014; see also <https://www.copyright.gov/reports/db4.pdf>). For example, data set tags used may in at least some instances be a sufficient enhancement to make the set protectible since the value point often lies with the expressive human input that lies with the tags to give life and meaning to the sample.

For more on general principles of copyright law applied to AI, see [Practice Note, Artificial Intelligence Key Legal Issues: Overview: Copyright](#).

Patent

Patent law is generally not a good fit to protect data. As a general rule under the “printed matter” doctrine, printed matter (or in this case specific data items) alone does not impart a patentable feature in a patent claim absent some functional relationship between the printed matter and the substrate (see, for example, *Praxair Distribution, Inc. v. Mallinckrodt Hospital Products IP Ltd.*, 890 F.3d 1024 (Fed. Cir. 2018)). This 18th century legal construct—that one cannot patent a book simply because the words are new—still applies today in the digital world. For example, a record or CD is not patentable simply because it has a new song on it. It is therefore similarly very challenging to seek patent protection for novel data items that make up a training set.

However, patent protection may be available when a novel and nonobvious process is developed to condition **how** a training set is developed. This could include, for example, a process that establishes parameters for inclusion or exclusion of data into the training set. For example, in the context of a system that takes readings from thousands of sensors, a novel and inventive system that identifies which sensor outputs are to be included in a training set and then tagged may potential result in meaningful patent protection. In such cases, the printed matter doctrine at least is unlikely to be an issue. Such protection, however, may not be helpful in the situation where the value lies in the original data as stored, without the need of further manipulation.

Organizations and their counsel will still need to consider other issues commonly found in protecting software like patent-eligibility, divided infringement, and adequate disclosure. For more on these principles applied to AI, see [Practice Note, Artificial Intelligence Key Legal Issues: Overview: Patent and AI Patent Infringement](#).

Trade Secret

Trade secret law is potentially available when the data training set is kept secret and the value of the data lies at least in part in the fact that it is secret. Often the process of monetization will require disclosure, however, which may not make keeping the training set secret a practical matter. Thus, trade secret law may also prove unhelpful in the protection and monetization of the AI/ML system’s operation, which will ultimately be public. For more on the requirements and challenges of protecting AI as trade secrets, see [Practice Note, Artificial Intelligence Key Legal Issues: Overview: Trade Secret](#).

Contract

Contract law is one of the best ways available today to protect and monetize AI/ML systems through licensing the right to use trade secrets. However, formulating a valuable contract requires that the matter being licensed is properly the subject of a license, and that an agreement is carefully formulated to be enforceable.

Other issues associated with owning and licensing data training sets (such as ensuring proper compliance with third party rights, open source agreements, regulatory limitations, and privacy laws) are beyond the scope of this discussion.

Technological Protection Mechanisms

In some instances, IP law may not protect AI sufficiently. Entities may instead use certain technology to protect and monetize AI data sets. For example, entities have used CAPTCHA, a program intended to distinguish human from machine input, for years to block bot crawlers from scraping websites. Entities and their counsel should therefore consider creative technological solutions where possible, which may in turn result in a separately patentable invention. In this way, rather than investing in expensive and potentially uncertain litigation, entities may use this technological alternative to block unwanted data mining that could compromise their AI IP.

Protecting the ML Algorithms or Models

In some instances, the ML algorithm or model the system implements may be an off-the-shelf algorithm or model that is akin to a word processor. In other cases, the software that implements the algorithm or model may be unique and innovative.

When the algorithm is innovative, as with other forms of software, entities and counsel should consider copyright and patents law to protect the algorithm, subject to the usual limits on such protection. Trade secret and contractual protection may also be available. For more, see [Practice Note, Legal Protection of Software](#).

Copyright

As noted above, copyright protects an expression, not an underlying idea. Therefore, while copyright may potentially protect the specific implementation of an algorithm or model as embodied in software code, nothing will stop other entities from writing their own code and resultantly implementing their own expression in a way that mimics an AI system. In particular, if the algorithm or model is off-the-shelf and published and

freely available, copyright law may offer limited protection unless a tremendous amount of programming acumen is involved to replicate the AI system.

Patent

An algorithm is considered an “abstract idea” under 35 U.S.C. § 101 in the current patent-eligibility paradigm. Therefore, if an entity pursues AI patent claims, it will need to claim and disclose more than merely the algorithm itself. The patent claims will also need to cover a practical application of the algorithm that is more often found as part of a larger system as discussed below. For more on AI patentability issues, see [Practice Note, Artificial Intelligence Key Legal Issues: Overview: Patent](#).

For a collection of resources addressing subject matter eligibility, see [Section 101 Patent Eligibility Toolkit](#).

Trade Secret

For trade secret protection, a key component will be keeping the algorithm secret, which may or may not be possible depending on the component’s use. Under some regulatory schemes and commercial situations, it may not be acceptable to keep the algorithm or model secret, such as where securities laws may require that the components of an index be publicly disclosed. Similarly, in order to obtain approval for a new medical device or diagnostic system, the FDA may require disclosure of how such device or diagnostic system operates.

Contract

Contract protection is often the best form of protecting ML algorithms, subject to the same caveats discussed above.

Other Technological Mechanisms

In some instances, other forms of intellectual property protection may simply be insufficient or too uncertain, and instead a technological solution, like requiring a user to enter an activation key to enable the software code, can be used. This is often the case when licensed with other software, like a commercial word processor or spreadsheet program. Technological solutions can avoid the need of costly and potentially uncertain litigation to limit the risk of misappropriation in the first instance.

Protecting the Input Query and Processes Used for Generation

Protecting the form of the input query may be the most challenging aspect of ML IP. As a practical matter, any third party using the system, such as a customer, will need to know

Protecting Artificial Intelligence Innovations as Intellectual Property: Opportunities and Pitfalls

how to formulate a query. The input query form will therefore not be secret, making trade secret protection unavailable.

An input query is also not likely protectible under copyright law because the process of entering the query is a functional requirement. Similarly, without something more than the query itself, it is unlikely to be protectable under patent law as an abstract idea under Section 101.

Instead, what might be potentially protectable under copyright law (as software) or patent law is any unique method developed to **generate and formulate** a query for the user to provide to the ML system. That is, to the extent value lies in the formulation steps, copyright or patent law may be useful to protect ML input queries.

Protecting the Output Result

Under current law, to the extent a ML system generates an output result and the computer is considered the “author” or “inventor,” U.S. copyright law and U.S. patent law do not allow for protection (see [Practice Note, Artificial Intelligence Key Legal Issues: Overview: AI IP Ownership](#)).

Contract law or trade secret law (and limiting access to the output result) are therefore the only forms of IP protection likely available.

Protecting the Larger ML System

To the extent possible, framing the innovation in terms of a larger system, including how the query is

formulated, the ML system itself (including the form of query, the algorithm or model), the training set and the form of output, and how the output is used to achieve a desired result, may be the best way to protect AI/ML innovation.

Larger systems incorporating the elements discussed above are potentially protectible under copyright and patent law, subject to the caveats discussed above. The system may also be subject to technological locks and contract protection as with any other software system. Trade secret protection may be potentially available for some components of the system that may be necessary to make the system work, but not otherwise necessary to disclose.

** Mr. Macedo is a partner at Amster, Rothstein & Ebenstein LLP where his practice focuses on all facets of intellectual property law, including patent, copyright, licensing computer implemented inventions such as machine learning and artificial intelligence innovation. He can be reached at cmacedo@arelaw.com. He would like to thank Herbert Blassengale for his assistance in preparing this article, and Daniel Dardani from MIT for his thoughts and input in developing the subject matter used to prepare this article.*

About Practical Law

Practical Law provides legal know-how that gives lawyers a better starting point. Our expert team of attorney editors creates and maintains thousands of up-to-date, practical resources across all major practice areas. We go beyond primary law and traditional legal research to give you the resources needed to practice more efficiently, improve client service and add more value.

If you are not currently a subscriber, we invite you to take a trial of our online services at legalsolutions.com/practical-law. For more information or to schedule training, call 1-800-733-2889 or e-mail referenceattorneys@tr.com.