

Professional Perspective

ChatGPT: IP, Cybersecurity & Other Legal Risks of Generative AI

Dr. Ilia Kolochenko, ImmuniWeb, and Gordon Platt,
Law Office of Gordon Platt, P.C.

**Bloomberg
Law**

[Read Professional Perspectives](#) | [Become a Contributor](#)

Reproduced with permission. Published February 2023. Copyright © 2023 Bloomberg Industry Group, Inc.
800.372.1033. For further use, please contact permissions@bloombergindustry.com.

ChatGPT: IP, Cybersecurity & Other Legal Risks of Generative AI

Contributed by [Dr. Iliia Kolochenko](#), ImmuniWeb, and [Gordon Platt](#), Law Office of Gordon Platt, P.C.

The globally publicized launch of ChatGPT has skyrocketed interest in intelligent chat bots, such as the flagship ChatGPT, and other forms of [generative AI](#). The ongoing media buzz has generated all kinds of overhyped and sensationalistic predictions around the alleged capabilities of contemporary AI, including some amusing ones, like the [proclaimed end of the legal profession](#).

One media prediction derived from the ChatGPT hype, however, seems to be perfectly correct: in 2023, venture capital and other investors will probably start pouring cash into AI startups, despite the unfolding recession in the technology sector. As a probable consequence—to the surprise of some enthusiastic “AI-will-replace-it-all” evangelists—law firms will likely have even more work than before, stemming from complex multijurisdictional litigation related to the development and use of generative AI.

This article is aimed at shedding some light on a set of interconnected cybersecurity, compliance and legal risks created by intelligent chat bots—such as ChatGPT—and similar forms of generative AI, as well as offering solutions on how to address them.

Training Data

Generative AI is about machine learning and machine learning is about data required for training. The trained AI model is only as good as the underlying data, both in terms of the quality and quantity of the training data. Not to diminish the pivotal importance of experienced data scientists, who actually orchestrate the entire AI training process and subsequent continuous improvement of the derived AI model, but without high-quality data even the most talented and experienced team of data scientists is poised to fail.

Depending on the underlying purposes of an AI system, different types and volumes of training data may be required. For instance, a technology assisted review (TAR) system, built to spot the most frequent problems with commercial contracts, will likely require a large set of diversified contracts for training in order to produce reliable and comprehensive results. Importantly, if the training data set does not include a certain type of contract or flatly misses a jurisdiction-specific clause, then the automated review of such contract will likely either overlook the problem completely or erroneously label a valid document as problematic. In sum, the more training data you have—both in terms of quality and quantity—the more reliable and efficient your AI system will be.

Foreseeably, a considerable number of AI startups and even tech giants may decide to engage in massive data scraping activities on the internet to outperform competition by collecting as much data as feasible. Data scraping is a fairly primitive data collection technique usually performed by automated software to browse and download various types of content including but not limited to various product databases, social media posts, news articles, books, scholarly articles and essays, court decisions and law texts, pictures, or even video and music. The data is commonly scraped from publicly accessible web resources, ranging from the Library of Congress and Wikipedia to smaller niche-specific websites, for instance, blogs or social networks of law firms with valuable analytical content originally produced by humans for humans.

Illustratively, according to the answer about its own training data, ChatGPT [claims](#) that its model “has been trained using a vast amount of text data, including books, articles, and other written material.” Eventually the scraped data will be used to train formidable AI models, which on their side purport to ultimately outperform and replace the content creators—at least in cost efficiency—even if tech companies say otherwise. The multidisciplinary stack of interconnected problems caused by scraping of training data is discussed below.

Intellectual Property

A heated debate now rages among US legal scholars and IP law professors about whether the unauthorized scraping and subsequent usage of copyrighted data amount to a copyright infringement. If the view of legal practitioners who see [copyright violations](#) in such practice prevails, users of such AI systems may also be liable for secondary infringement and potentially face legal ramifications. Japan, interestingly, has taken the lead in amending its copyright laws which have

weakened the protection of creative content authors when their works are used for AI training purposes. The EU copyright legislation currently provides some narrow exceptions allowing data mining of copyrighted materials for scientific research purposes including AI training, however, it is [contended](#) that those exceptions rather hinder data mining and provide copyright owners with additional rights and powers.

The authors of this article believe that, in most countries, the existing copyright legislation will unlikely provide a solid protection for authors of misappropriated content exploited within the context of AI training purposes as long as the content is not reproduced directly in the derived work product. Of note, in some under-researched or novel areas of knowledge, where available data is scanty, the output of generative AI will likely resemble the original content, thereby increasing a plaintiff's chances of [prevailing in court](#) for copyright infringement.

Owners of creative content, however, have other legal instruments to protect their intellectual property from uninvited bots. The legal landmine—vastly ignored by unwitting AI companies that operate online bots for data scraping—is hidden in Terms and Conditions commonly available on public websites of all types. In contrast to the currently unsettled IP law and the copyright infringement dilemma, a website's Terms and Conditions are backed by well-established contract law and usually can be enforced in court relying on sufficient number of precedents.

Historically, a significant number of Terms and Conditions contain a prohibition of data scraping activities, having its roots in late 1990s. At the dawn of the dotcom epoch, unfair market players in the booming ecommerce sector aggressively scraped pricing and other product data from websites of their competitors to offer symbolic discounts and get higher rankings in then-nascent sales aggregators. In recent years, data scraping has become a persistent headache for ecommerce and other industries, as unwarranted usage of data may still bring a potent competitive advantage to unethical market players.

In addition to unscrupulous businesses, cybercriminals also regularly launch data scraping campaigns, for example, to create a copy of a well-known website and selectively infect its pages with malware prior to promoting the fake site via Google Ads. Consequently, most boilerplate Terms and Conditions for websites—abundantly available in free access—contain a clause prohibiting automated data scraping. Ironically, such freely available templates have possibly been used for ChatGPT training. Therefore, content owners may wish to review their Terms and Conditions and insert a separate clause flatly prohibiting all usage of any content from the websites for AI training or any related purposes, whether collected manually or automatically, without a prior written permission of the website owner.

On social networks, some ingenious technology experts argue that data scraping cannot be reliably attributed to a specific market player and thus can be safely used by innovative AI startups to advance technological progress. Omitting the ethical ingredient of such questionable assertions, it is important to note that large-scale data scraping campaigns can be reliably attributed to the AI companies that actually pull the strings, even if the latter insidiously exploit subsidiaries from offshore jurisdictions or rent volatile cloud infrastructure as proxies.

Therefore, inserting an enforceable liquidated damages provision for each violation of the no-scraping clause, enhanced with an injunction-without-bond provision, can be a tenable solution for those authors of creative content who are not keen to provide the fruits of their intellectual labor for AI training purposes without being paid for it or, at least, given a proper credit for their work. From a technical perspective, website owners may also consider implementing a bot-protection system, which are usually included into all modern web application firewall (WAF) offerings. Similarly, software engineers intent on protecting their open-sourced software from being used for AI training purposes, may insert the foregoing provisions into the software license agreements before posting their source code on GitHub.

In addition to Terms and Conditions, another legal trap for companies that run internet-wide data scraping campaigns lurks in unfair competition legislation. For instance, an AI startup that relies on clandestine data scraping techniques, while its competitors pay for training data, may likely find itself in crosshairs of market regulators and watchdogs of fair-trade practices. Moreover, in some countries, a company that knowingly relies on a generative-AI service provider—being fully informed of or willfully blind to the unauthorized data scraping activities of the AI provider—to “smartly” outperform its competitors, may also face serious sanctions or even criminal prosecution for unfair competition.

Additionally, the end users of generative AI, whether a copywriting company or a consulting boutique offering creation of internal cybersecurity policies, may well be found in breach of contract if instead of producing—and actually billing for—qualified human labor, outsources the writing process to an AI system. That being said, there is nothing wrong per se when

busy human experts leverage generative AI for acceleration and intelligent automation of their work, but such practices should be transparently disclosed to the end customer for well-informed decision making.

Proposed Regulatory Response

Whilst the authors of this article firmly believe that continuous innovation with the development of AI technologies is the future for virtually all industries, they are likewise convinced that innovation should be adequately regulated to produce socially desirable outcomes, sustainability and social equity. Otherwise, AI companies may swiftly gain unfair advantage by monopolizing existing human knowledge –without paying a dime to its creators–to eventually sell the knowledge back to humanity.

To comprehensively address the challenge, lawmakers should consider not just modernizing the existing copyright legislation, but also implementing a set of AI-specific laws and regulations. That being said, to avoid counterproductive overregulation that will eventually stifle the necessary innovation, it will be useful to invite experts and entrepreneurs from all concerned industries to make comments, contributions and proposals on the AI legislation to ensure a fair balance of everyone's legitimate interests.

Regulation of AI may have two primary prongs: mandatory disclosure of data sources used for AI training and creation of a new fundamental right for creators of original content in addition to the contemporary notion of copyright protection.

Transparency of data, used for AI training, is not only of extreme importance for better predictability and explainability of AI systems, but also for the prevention of unwarranted or prohibited data scraping. Therefore, AI vendors of a certain size may be required to disclose all sources of their training data on a regular basis, ensuring auditability and accountability.

Furthermore, individuals and legal entities may be provided with an additional, explicit right to restrict the processing of their creative content for AI training purposes without first conveying an informed consent. This right may be similar to protection of personal information, where a prerequisite permission is necessary to collect and process personal data. Before May 2018, countless opponents of EU's GDPR energetically argued that the new legislation would stifle innovation and cause a pan-European recession, but in reality, both data subjects and data controllers now operate in a healthier, safer and more equitable environment.

Remarkably, being inspired by the proven success of GDPR, many countries already have or are currently implementing similar privacy legislation. Creation of the new right, will likely incentivize many professionals, such as lawyers, doctors or journalists, to continue publicly sharing their invaluable knowledge and unique expertise with other humans, as well as to sell or license it to AI giants, supporting equitable innovation and technical progress.

Cybersecurity and Compliance Risks

Assuming that lawmakers introduce the necessary AI regulation in the near future, the generative AI industry will undoubtedly flourish and bring sustainable value to society. The inevitable success will, however, come with a set of interconnected cybersecurity and compliance problems that require special attention and advance planning to be properly mitigated.

First, AI providers will likely become a desired target for nation-state cyber threat actors and sophisticated cyber mercenaries. Intercepting company- or user-specific input to an intelligent chat bot can disclose a plethora of confidential information, ranging from personal data and trade secrets to law enforcement intelligence and classified governmental information. On top of that, sophisticated intruders may go one step further by providing modified and backdoored answers to specific users or organizations, aimed to cause damage. Despite that supply chain attacks are far from being novel, with the rapid proliferation of AI, third-party risk management (TPRM) programs will become even more critical.

Second, externally available training data can be poisoned with deliberately incorrect and misleading information by hackers, organized crime, or even sovereign states in their proliferation of international disinformation campaigns to, for example, interfere with democratic elections. Many internet archives and public libraries do not have budgets to protect their numerous web applications, let alone to prevent intrusions into their corporate networks. As a result, gaining access to their web systems is much easier and faster than breaking into an e-banking application or an average e-commerce platform. Cybercriminals may stealthily modify some data, perfidiously altering crucial facts, vital knowledge, or critical weights and numbers.

As a consequence, an intelligent chat may tell a patient to take a fatal dose of prescribed drug, instruct a database administrator to delete data from a state election database instead of repairing it, or distort historical facts about atrocities of war crimes prior to answering a question from college students. Whilst such attacks require advanced technical skills and significant financial resources to be operationalized with a sufficient scale and necessary frequency, their consequences can be truly catastrophic.

Third, unless you own and host your own version of ChatGPT, whenever your employees enter anything into an intelligent chat, you should carefully assess your legal risks. How will your data be used and who will have access to it? Will it be shared with any third parties? This is a pivotal question given that when using an intelligent chat your software engineers may inadvertently copy-paste your proprietary source code, your corporate HR may unconsciously enter sensitive personal data, and your C-level executives may unsuspectingly disclose your upcoming acquisition transaction that may significantly impact your stock price.

Therefore, to protect your trade secrets, regulated data, and other confidential information, you should consider establishing and promulgating an internal policy on permitted usage of chat bots and similar generative AI systems. Regular training of employees will help to increase their awareness about the risks to facilitate compliance.

Finally, generative AI may perfidiously hide in your office even if you never used or considered using such system. A rapidly growing number of SaaS providers—including tech giants—now insert a clause into their Terms of Service that expressly authorizes them to utilize your corporate data for AI training purposes. For instance, when you upload a contract for signature via a SaaS document management platform or send an email via your cloud mail provider—its entire content may be silently used in way you never expected.

Conclusion

Generative AI, such as ChatGPT, can bring immense value to the humanity by intelligently automating countless routine but time-consuming tasks and processes, eventually enabling people to focus on more important things that truly deserve human ingenuity. The proliferation of AI must, however, be properly regulated to prevent tech monopolies and to ensure sustainable and fair development of society. Before lawmakers enact the necessary legislation, authors of creative content can consider using simple but efficient techniques described in this article to protect their intellectual property from misappropriation by some AI companies.