

## Contracting for Generative AI

March 2, 2023

### **Announcer**

Welcome to Mayer Brown's Tech Talks Podcast. Each podcast is designed to provide insights on legal issues relating to Technology & IP Transactions, and keep you up to date on the latest trends in growth & innovation, digital transformation, IP & data monetization and operational improvement by drawing on the perspectives of practitioners who have executed technology and IP transactions around the world. You can subscribe to the show on all major podcasting platforms. We hope you enjoy the program.

### **Julian Dibbell:**

Hello and welcome to Tech Talks. Our topic today is a hot one: generative AI and the contracting and other legal issues that can arise for users and data providers with respect to this potentially disruptive technology that everyone's talking about.

I'm your host, Julian Dibbell. I am a senior associate in Mayer Brown's Technology & IP Transactions practice. I'm joined today by Marina Aronchik and Richard Assmus. Marina is a partner in the Technology & IP Transactions practice as well. Marina advises leading companies on a wide range of regulated and non-regulated industries on cutting-edge technology issues and transactions. Rich is a member of the Intellectual Property, Brand Management and Litigation practice, and he co-leads the firm's Technology & IP Transactions practice. Rich has a balanced intellectual property transactional and litigation practice, and he has been advising clients for several years on emerging issues in AI.

So, let me ask you, Rich, we've been focusing on the impact of AI for a few years now, us in the technology law area, but the topic of "generative AI" has suddenly catapulted to the center of discussions around AI. So, what is "generative AI" and why has it garnered so much recent attention.

### **Rich Assmus:**

Thanks Julian and hello everyone. The idea that AI tools might generate output has been around for a while, but the early commercial applications of AI were really more focused on decision tools – using AI, for example, to assist with insurance underwriting, or producing relatively simple output (like answers to straightforward questions). Those AI tools certainly generated output of some kind – for example, a recommendation denying or accepting insurance coverage, but not in the sense of output that might substitute for the creative content that a human might author him or herself.

So, when people talk about "generative AI", they mean AI that creates some type of consumable content, whether that content is computer code, text responses, images, or even music.

Probably the number one thing that has created so much buzz around generative AI recently is the latest iteration of ChatGPT. That's a language model that uses AI to deliver an uncanny chat experience that is very good at answering questions with well-crafted natural language answers. ChatGPT can even tell short stories if you give it the right prompt, or even write poems. But ChatGPT is by no means the only generative AI tool out there, and as we'll talk about later, other tools like GitHub Copilot and certain image tools have made the news, and drawn the ire of intellectual property rightsholders. Those rightsholders are claiming various violations in the way those tools are trained and in the content of their output.

**Marina Aronchik:**

And, if I may jump in here, Julian, there is another aspect to generative AI that, frankly, I did not appreciate until very recently, which is that some of the generative AI tools use the language model that not only answers questions but also asks questions and makes requests of the user, much in the same way that you'd see in a human conversation. There are transcripts of conversations with some AI-based chatbots where the chatbot – and I'm paraphrasing here – responds with something along the lines of *how would this make you feel? What would you do? Or This question makes me uncomfortable – let's switch topics.*

Being an analytical person, my first reaction to being surprised at this was really to ask myself, why was this surprising? And the reason that I think that it was surprising was that for a number of years we were really analyzing AI and thinking about AI and working on contracts for the commercial applications – we always thought of AI as responding or processing training data or production data and providing output, without this loop back of AI output basically becoming the input to the user. And of course this makes sense if you think of AI using natural language, and a natural part of a human conversation is for humans to ask one another questions either to seek information or to maintain the flow of a conversation. But again when Rich mentioned the uncanny experience, I think that's really part of it.

**Julian Dibbell:**

I totally agree, Marina, I think you have hit on a really important point that this is a new kind of AI that's really taking things to a whole other level, a much more complex level potentially. And obviously, this is why it's getting so much attention, why it's in the news. How is it impacting our clients? How is it impacting businesses that may want to use this technology?

**Marina Aronchik:**

That's a great question. We could spend hours talking about this, but I'll take a shot at it. So, let me start by saying this: I came across a line in the New York Times yesterday about AI that I think is key to this discussion, and the idea is that AI technology that we thought was going to arrive in 2033 actually arrived in 2023, meaning today. So, what we've been seeing in the last few years are relatively limited applications of AI and the output of those tools was really in the nature of the data the company might provide to itself, rather than using pre-trained algorithms. So, based on a variety of factors and the profits that are at stake there, there was really room, I think, for caution in adopting these tools, and thinking through the issues, which takes time.

So, the AI gold rush that arrived in 2023 is very different. These tools are creating content that, as Rich mentioned and I'll reiterate here, substitutes for content that a human might create himself or herself, and also these tools are broadly available and, importantly, they represent immense financial opportunity and tremendous investment. So in part because of these factors and a number of other factors, we're now at a point where there's fast and widescale adoption of the tools and that things are moving so quickly that there is really not a choice to apply the caution that we'd typically see or that we have been seeing in the prior years. The game now is to try to keep up and adjust and account for risks as they come up. It's sort of like changing tires on a moving car.

**Julian Dibbell:**

Right, and so we've got business users at these organizations jumping in and using this technology whether or not the organization itself has had a chance, and who has had a chance, this is coming so fast, to really sit down and consider and establish an overarching policy that their users understand for using this, right? Obviously, it's time for them to start thinking about the risk, and I presume they are starting to think about these risks, right, Rich?

**Rich Assmus:**

Yes, they are. We're getting calls on this with greater and greater frequency, and clients are really considering a number of risks here and today we want to focus on two distinct legal issues that arise, specifically under copyright law. To be clear, there are other issues as well. Copyright protects a broad array of creative works, including computer code, images, and music, and virtually all of the input and output of the generative AI that's out there is within the ambit of copyright law generally.

And really the first legal issue that arises is about the training data that is used to refine the AI tool. Without that training data, an AI tool really isn't useful. To take ChatGPT as an example, it's so good at providing these natural language responses because it's been trained on a huge database of actual human language. And similarly, an AI coding tool like CoPilot, it's so good at providing working code because it's been trained on working code created by humans. Many AIs are essentially mimics and they need something to mimic. All of these generative AI tools have the same feature – they work as well as they do based in part on the quality of their training data.

So really the question is whether the AI tool has been trained on input that's subject to third-party rights, and whether that training use is an infringement of any of those rights. And to be clear, we're assuming here that the AI provider did not license the right to use the training data from any of those rightsholders. That would always be an option for a tool provider and would set this particular risk aside.

**Julian Dibbell:**

Okay, but help me understanding something here. Why is infringement an issue in the first place? My understanding is that these new AI tools are generally trained on information that is publicly available, right? Repositories of text like Wikipedia, for example, or images or code that anyone can go and look at or even, in the case of open-source software, make use of. What is different here?

**Rich Assmus:**

Right, so that idea that if it's publicly available, that means I can use it, that's one of the myths I think that drives copyright lawyers nuts. Just because something's available on the web or on some other form doesn't necessarily mean you have the right to use it for any purpose. Sometimes those websites that you might be visiting have terms of service that arguably bind you to the way you can use that material. In other cases, it would just be beyond the scope even of any implied license you have. So, just because something's available on the web doesn't mean that there's not a copyright owner behind that information or text or whatever it is. For example, you can probably stream any song you can think of on Youtube, but that does not mean that you would have the right to use that music in any way. So that's really why AI tool providers need to think about the rights they have to use their training data, even if that training data is available at the touch of a computer keyboard or mouse.

It's probably worth here going back to first principles of copyright law. Assuming a given work is protected (and let's assume that here that the training data is so protected, it's a fairly safe assumption), a copyright owner has various exclusive rights, including the right to make copies of a work. Given the way that computers operate, it's a safe assumption also that each work within the training data set is copied during the course of the training. For example, that training data might be copied from the server on which it's hosted to the permanent storage of the computer that's doing the training, and then from there onto the working memory of the computer. So we're in this situation where the use of copyrighted training data at a minimum is going to implicate the rights of a copyright owner. That's by no means the end of the story, but it's the basic reason why we're talking about this risk at all.

**Julian Dibbell:**

But there is an exception for fair use, right?

**Rich Assmus:**

That's right, there's a very important statutory exception called fair use and it makes certain uses of copyrighted works not infringements at all. There's a four-part test for what constitutes fair use, and the statute and the case law really center around what amount to policy choices about not allowing the very broad brush of copyright law from stifling creativity.

I'll quickly run over those factors for you. The first is the purpose and character of the use, including whether such use is of a commercial nature or is for nonprofit educational purposes. With respect to an AI tool that's in commercial operation, arguably the use of that training data is commercial in nature, although I don't think that's entirely clear.

The second factor is the nature of the copyrighted work and that really means whether it's something that is at the heart of copyright law, like a novel or music, or something that's less creative than one of those things. The third factor is the amount and substantiality of the portion used in relation to the copyrighted work as a whole. So, for example, in writing an article, you might quote a section of another article, and that would largely be considered fair use. One complication here is that the training data is typically used in whole, not in part.

And then, finally, the fourth factor is the effect of the use upon the potential market for or the value of the copyrighted work. I think that's an interesting question in the case of AI, where you have AI tools that are trained on code, for example, that are then creating code that could indeed substitute for the work of a human programmer.

So the question as it applies here, is whether that training data use should be deemed as fair use. We can run through the statutory factors, obviously each case is going to be a little different. We can make a judgment on how the courts would apply the factors, but right now we really don't have much guidance at all from the courts on this specific point as it relates to the use of training data as an input to AI.

**Julian Dibbell:**

You've painted a very lucid picture, Rich, of the potential copyright risks involved in the use of the training data that input into the AIs and we'll talk a little bit later about the output concept. Marina, I'm wondering, from your perspective as a contract lawyer, what does this mean for organizations that are trying to potentially do deals for the incorporation of these tools?

**Marina Aronchik:**

It means that while logically, it makes all the sense in the world, for the provider of the AI tool to take on the risk that, basically, the tool's training data or the use of the data was infringing, in practice, it's difficult for providers to quantify this risk that's an unknown that Rich just described. There are somewhat different issues that we also hinted at earlier depending on whether the specific training data that was used for the tool was publicly available data. Again there might be licenses. There almost certainly are licenses and limitations to the rights associated with that open-source licenses being at the top of that list versus those where the AI provider relies on commercial contracts with third parties for standalone data licenses, or data licenses that in effect are embedded in other product or service agreements (including cloud agreements). So, it's quite a thorny and complex issue.

And going back to the contract governing use of the AI tool, the first question would be, well, what contract applies? Do you even have an opportunity to negotiate that contract, or are you bound by the online terms of use that obviously for companies that we're advising that depends on where that tool was implemented and how it's used. And if you do have an opportunity to negotiate the relevant license, depending on where you come out in negotiations, there is a real risk that IP infringement or other third-party claims arise simply based on your use of the AI tool, and that might be the risk that a business needs to take on, meaning you as a customer using that AI tool are now taking on the risk of infringement and again that goes back to the contract and whether or not the provider of the tool is willing to take on the risk. Anticipating a business point that will probably come up – this isn't this the same risk that we have with other software. In most negotiated and sophisticated agreements, the provider of the software will generally give you a non-infringement rep and warranty and there's going to be an indemnity associated with that and that's in part because these providers have comfort that they can understand their risk and they can manage that risk. That's not at all what we're seeing with AI tools.

**Julian Dibbell:**

Right, and to be clear, the uncertainty is not just based on the fact that there's such a vast amount of data

that it's really hard to track what the license rights might be, but it's also that the training that inputting of data is ongoing, right, it's continuing while the tool is used.

**Rich Assmus:**

That's right. One of the things that means is that clients that have data need to be worrying whether or not their business partners are using any of their data for training purposes, and whether their agreements address that. As an example, many form contracts have the concept that you could use the data for internal business purposes, it's a very common phrase you see in agreements, and one of your counterparties might say, well, the training of our AI tool that we're going to be marketing to your competitors, that wasn't internal business purposes. We're not giving that data to any third party. So we're in a situation where many form contracts really aren't cognizant of the AI risks simply because the technology is so new.

**Marina Aronchik:**

Right, and I'd add to that it's not just the training data or even production data that we're providing. When you think about chatbots and this is true for other AI as well, these tools are likely soliciting additional information from the client above and beyond the training data and the production data, in part through these, what we view as natural conversations. I rely on an analogy here so I would think of it as a combination of a cloud tool (for the core functionality), with a human professional service (training, maintenance and support or otherwise), where you could well have a conversation with a person soliciting this additional information and prompting further input and potentially inducing you to take certain actions but the data that's being provided to the tool, we need to think of it in this full continuum, the full spectrum of interactions and uses – which depend both on the tool, and also the evolution of the tool over time.

And even if a company is comfortable with its own data being used for ongoing evolution of the tool, right, so we take a common scenario, here's my AI tool, I'm going to provide some data to it and of course I expect it to evolve over time, and that's acceptable to me for the particular business case, to the extent that a company is using third-party data (customer data or maybe even data of other vendors), that type of downstream use by the AI tool may not be permissible and almost certainly was not contemplated by the relevant agreement, so now you could be in a situation where there's potential breach of the obligation to these third parties and to the covenants that you might have made to the AI tool provider by your right to use third-party data. So this all becomes very complicated very quickly and I'd add here the issues aren't limited to IP. There are potential regulatory concerns here and this is where we need to think of it on an industry-by-industry basis, but there could well be competition concerns or some of the regulators out there are worried about scenarios where Company A is providing production data, training data into the tool that same tool is used by a competitor in another industry, and so you could see how a tool could now be used to indirectly exchange information by training an AI tool in the same way both if you say the same information to another third party, an intermediary, that could be a competition concern. And, again, that's just one of the examples of the challenges we should be thinking about.

**Julian Dibbell:**

I could see a reasonable business coming to these tools and saying look, the conservative, the prudent approach is to, yes, we want to use the tool, but we are opting out of the part of the agreement where it says you, the tool provider, can use our data, whether it's through the connected systems or like you were saying, just solicited input, right?

**Marina Aronchik:**

Right, and it's a nice option and it could be successful particularly with smaller providers who have more flexibility. For a larger provider, and again keeping in mind the nature of the tool, there may not be the technical flexibility to do so, it might be outside of the business model and it may be contrary to this idea that the AI tool is always evolving.

**Rich Assmus:**

And I think you've both sort of touched on this, there's an important question to consider here, really about trade secrets and whether or not either the data you're providing to the tool or the data the tool is soliciting is covering information that a company might view as a trade secret and whether a tool that's trained on trade secret data might be spitting out trade secrets at the back end in connection with the services it's offering. There's a lot of questions about the way the AI tool can act as an intermediary for disclosures that a company would never make in the ordinary course.

**Julian Dibbell:**

Right, and I just want to make sure our listeners understand because they may be familiar with more traditional kinds of AI where you're just submitting very kind of rote questions or queries and getting kind of constrained answers. As Marina was saying, it's a wide-ranging conversation. You're having people taking prompts, the machines are taking prompts from users to create images or software code, and those prompts can be freewheeling and they can reveal a lot about an organization so the idea is that a lot of confidential, important, and potentially trade secret information could be divulged. But that's probably yet another topic for a whole other episode.

Let me turn now to that other point that Rich alluded to and we talked about. There are two legal issues arising from generative AI, and the first relates to the input, the training data and we've talked about that I think pretty thoroughly. What is the other legal issue you were hinting at?

**Rich Assmus:**

The other issue relates to the output of the generative AI tool. You can think of an AI as an input/output machine. I think it's a little more complicated than that in that it's a bit of a loop. I was just reading an article today about these machines that have been built to play chess or to play the game of "Go" which is actually in some ways a more complex game than chess in terms of the number of potential board positions, and some of those Go-playing machines were actually trained by playing themselves, so there's definitely a circular loop here that makes the input/output model a little too simplistic, but I do think it's the right way to think about these copyright issues. And so really, what are the copyright issues about the output of the AI tool. There's really I think a number of subsidiary points. The first one is getting a lot of attention, which is whether the output of generative AI is itself copyrightable. Personally I don't think that

issue is going to have much practical importance, although it does point out an important gap in current copyright law, namely that copyright law, at least in the United States and much of the world, is designed to protect human authorship, and right now we don't know whether the output of generative AI really has a human author to claim rights. So, for example, if you gave ChatGPT a prompt, are you in some sense the author of the output? That's not clear at all. Without such a human author, those works are essentially public domain from a copyright perspective. Now to be clear, that doesn't mean they couldn't be treated as trade secrets, but you wouldn't be able to assert copyright rights in that output.

There's a few reasons I don't really think that's going to be too important, at least right now. First, a lot of the uses of generative AI will be integrated with expressions that are created by humans. For example, to take the coding situation, a human coder might use ChatGPT to fill in parts of the functionality of their program, or Github Copilot, but the code as a whole is going to be subject to human authorship and therefore protectable under copyright law. That protection is really very likely sufficient to prevent the uses a rightsowner will really care about. Similarly, stock images generated from a tool like StableDiffusion, they're likely to become part of a larger advertising piece, and those users may not really care about copyrightability of individual elements of their overall expression.

I think really the more interesting issue is whether the output of generative AI might be an infringement of one or more of the works in the AI tool training data set. And like the issue of fair use, we really don't have good guidance right now about how courts are going to consider that question. There however are several cases pending in the courts that have the potential to shed some light on these questions, including lawsuits pending against Github Copilot as well as the StableDiffusion image tool.

Just like the issues with AI input, these legal risks on output, they are prompting clients to re-visit their contracting policies with respect to generative AI tools, and to consider more generally what policies to put in place governing the use of these new tools.

**Julian Dibbell:**

While we await guidance from the courts, and that's going to be a bit of time coming I assume. Marina, what can we do in the meantime with contracts? Are there contractual solutions to this output issue?

**Marina Aronchik:**

Well, the first thing to note on this is one must serve with due diligence and understanding what terms are already in place, because as we mentioned earlier, some of the tools might already be in use or may have been used in the last twelve to eighteen months, and so it's important to understand what terms apply to those tools such that you can then determine how to treat output and exactly what those risks may be, and then you can think about remediation.

Going forward, with a caveat that this is a rapidly evolving area with intricate questions, I think the solution is to at least allocate the rights contractually, to the extent that you can do so through negotiations with respect to the output as between the parties. This does not solve the challenge of the potentially infringing nature of that output, but customers should be negotiating, again to the extent that you can, indemnities to address this risk, but at least allocating the rights as between the AI tool provider and the

customer is an important step. And, maybe even more importantly, any outbound licenses and downstream commitments that you're making, that a company is making, to your own customers, to other users, to other third parties regarding this output or regarding the products or tools or deliverables that incorporate the output generated by AI should probably be subject to this uncertainty regarding IP rights and the potential infringing nature of these components. Perhaps here, again looking to an analogy, we look to open source exceptions where generally the use of open source remains subject to the relevant licenses and there is not an expectation that a company using open source make its own reps and warranties regarding open source. So here you could be including disclaimers as to the output of the AI tools, but importantly, one must understand what these components are.

**Julian Dibbell:**

Understanding all of this is a big job for all of us in the months and years ahead. Thank you for getting us started on that process, Rich and Marina. We've clearly only touched on some of the IP and contracting issues relating generative AI. There are all kinds of other issues potentially on the horizon. I really look forward to continuing this conversation with you all in the future on coming podcast episodes, and hearing from our listeners about this.

**Marina Aronchik:**

Thanks for having us on, Julian.

**Rich Assmus:**

Thanks.

**Julian Dibbell:**

Alright. Well, listeners, if you have any questions about today's episode, I mean it, I would love to hear from you, or if you have an idea for an episode you'd like to hear about anything else related to technology and IP transactions and the law – please email us at [techtransactions@mayerbrown.com](mailto:techtransactions@mayerbrown.com). Thanks again for listening.

**Announcer**

We hope you enjoyed this program. You can subscribe on all major podcasting platforms. To learn about other Mayer Brown audio programming, visit [mayerbrown.com/podcasts](http://mayerbrown.com/podcasts). Thanks for listening.