



# Is ChatGPT just a copycat?

As ChatGPT turns one year old, there are growing questions about the way it draws upon creative works to compete with the authors of those very same works.



**Simon Chesterman**

Artificial intelligence (AI) has always depended on access to data.

Today's large language models (LLMs) are trained on, essentially, the entire Internet. Much of that is public domain material outside the realm of copyright. It also includes pirated works that should not be there and material that was shared to be read but not copied.

The scale of these models transformed public debate about the impact of AI with the release of ChatGPT by OpenAI in November 2022, quickly followed by competitors such as Google's Bard, Anthropic's Claude and Meta's Llama.

Excitement and trepidation abounded at the ability of these systems to respond to natural language queries with human-like responses – in text as well as images. Goldman Sachs breathlessly reported that generative AI could increase global gross domestic product by a whopping 7 per cent.

As Dr Mark Cenite of Nanyang Technological University argued last week ("ChatGPT can do more things than we realise. There's just one problem", Oct 17), access to copyrighted materials will significantly improve the quality of what AI can do.

But how (if at all) should the authors and artists whose text and images trained the underlying models be recognised and compensated?

#### A PIRATE'S LIFE

The use of pirated or illegally obtained material appears to be a simple case of theft. However, this has been notoriously difficult to prove.

A lawsuit brought by Getty Images against Stability AI, for example, includes images with distortions of the watermark Getty uses to protect its product.

Given the secretive nature of much model training, proving infringement is rarely as easy as this. Even in the Getty case, infringement may need to be proven on a case-by-case basis, establishing substantial similarity image by image – rather than the systemic copying of 12 million pictures alleged by Getty.

Even if infringement can be established, fair use is a defence that balances the rights of creators and the interests of the wider public in distributing and using their works.

It generally considers the purpose of the use, the nature of the work, the amount used and the effect on the market for the original work. When an individual records a televised broadcast to watch at a later time, for example, that can be considered fair use. Projecting such a recording for an audience and charging for tickets would not be.

A key question is whether using data to train models, which then produce works that may compete with the creators of those data, constitutes fair use.

When Google began scanning vast quantities of books in 2002, there were challenges that this infringed copyright. Google was, for the most part, successful in arguing that it made snippets of the information available but was not itself threatening the market for the original works.

By contrast, the ability of generative AI to produce text and images that do compete directly with past and present works is central to several other lawsuits currently under way – involving prominent authors such as John Grisham, Jonathan Franzen and Elin Hilderbrand who are suing OpenAI, the creator of ChatGPT.

#### DATA MINING FOR GOLD

Singapore is an example of a jurisdiction that has tried to thread this needle through legislation.

Amendments to the copyright law in 2021 include a new permitted use to make a copy of a work for the purpose of "computational data analysis". That includes extracting and analysing information and using it to "improve the functioning of a computer program in relation to that type of information or data".

Lawful access to the underlying data is still required, but it appears more open to data mining and model training than traditional conceptions of fair use.

It is also wider than the "non-commercial" text and data analysis exception adopted in Britain in 2014, or the "text and data mining" exception adopted

by the European Union in 2019. An information sheet produced by the Intellectual Property Office of Singapore explicitly states that the provision is intended to allow "training machine learning". Yet, analysing text or images for the purpose of making recommendations or optimising workflows is quite distinct from using those texts and images to generate more text and images.

Whether an author could successfully bring an action against OpenAI or another model developer would depend on proving infringement and that no exception applies.

Given the secrecy around model training, the first hurdle could be a challenge – though scholars such as Mr Peter Schoppert have helpfully created tools that enable authors to find out whether their works are in the databases used for model training.

If those works were used without lawful access, infringement seems clear.

In other cases, the Singapore exception for data mining is broad. Nonetheless, the law specifies that the materials should not be used for any purpose other than computational data analysis. If they are used to create new artistic works that compete with the original works, that may fail to satisfy the fair use test too.

This is easier to show in the case of visual arts.

When I met students at Singapore Polytechnic for an "AI Manga" competition in September, they were excited about the AI tools they could use to enhance their work – and concerned that no one would be willing to pay for it, because it is now so easy to create your own content.

For text, it may require proving a similar economic impact. Using AI to draft an e-mail or complete an assignment might not meaningfully dilute the sales of books by Rachel Heng or Kevin Kwan. (Although Crazy Rich Ais is a title I could get behind.)

It may be clearer in the case of journalists, whose daily scribbles are vacuumed up and repurposed in a manner that has challenged the viability of newspapers around the world.

More generally, even if the actual damage is small, the principle of limiting the use of creative works in a manner that does not allow AI to decimate the market for those works is worth fighting for.

#### A NAPSTER MOMENT?

This is no longer a hypothetical problem. In addition to diluting

**When Google began scanning vast quantities of books in 2002, there were challenges that this infringed copyright. Google was, for the most part, successful in arguing that it made snippets of the information available but was not itself threatening the market for the original works. By contrast, the ability of generative AI to produce text and images that do compete directly with past and present works is central to several other lawsuits currently under way – involving prominent authors such as John Grisham, Jonathan Franzen and Elin Hilderbrand, who are suing OpenAI, the creator of ChatGPT.**

the value of human authors' works, it is possible that they will simply be swamped by the volume of generative AI produced competitors.

Amazon, today one of the world's largest publishers of books, became so overwhelmed by submissions that it imposed a limit that its self-published authors may now publish "only" three books a day.

So, what happens next? The music industry offers interesting parallels.

It also went through a period of unrestrained piracy in the early digital era, which radically transformed the economics of copying and gave rise to file-sharing services such as Napster.

Lawsuits and legislative changes led to most media platforms adopting copyright policies and takedown protocols, while those like Napster were

shuttered completely.

It is possible that a similar evolution will take place in generative AI.

Adobe, for example, has built its Firefly tools using training sets consisting only of public domain and licensed works. Shutterstock has also committed to building AI tools with a Contributor Fund to compensate artists.

Other approaches are possible, such as the way YouTube allows certain usages of music and other copyrighted material by sharing advertising revenue with owners of the original work through its Content ID system.

Another option is provision for content creators to "opt out" of being scraped for their data, either through the site's robots.txt file or registering its Internet Protocol address.

#### GOOD MODELS BORROW, GREAT MODELS STEAL

T.S. Eliot once observed that "good writers borrow, great writers steal".

Eliot was not, of course, condoning plagiarism. His larger point was to challenge naive idealisation of the creative process: in arts, as much as in science, each new thinker and writer builds on the work of those who have come before. Painters inspire and echo one another; writers offer variations on plots and structures that can be mapped and catalogued.

This is clearest in music, where the limits of the heptatonic scale and chord progressions mean that melodies will inevitably echo one another, as Ed Sheeran successfully argued in a case concerning similarities between his hit song Thinking Out Loud and Marvin Gaye's Let's Get It On.

In any case, it may seem pointless to argue that AI models should pay for the use of data when the entire Internet has already been absorbed.

In addition to the market for "legitimate" models, however, there is evidence that further refinement of those models and the training of new ones depends not just on the volume of data but its quality.

Early suggestions that LLMs might continue improving based on synthetic data that they themselves create have foundered on projections that such AI-generated data will "poison" future models. Presuming there is an ongoing market for data and the political will to regulate it, the idea that generative AI will have its own "Napster moment" is at least plausible.

So, before we celebrate ChatGPT's birthday, let us consider whether it might be fair for it to help pay for the party.

• Simon Chesterman is vice-provost (educational innovation) at the National University of Singapore and dean of NUS College, as well as senior director of AI governance at AI Singapore. His latest book is the novel *Artifice*.

The scale of large language models transformed public debate about the impact of AI with the release of ChatGPT by OpenAI in November 2022, quickly followed by competitors such as Google's Bard, Anthropic's Claude and Meta's Llama. PHOTO: REUTERS