



AI in the Publishing Industry
A Brief Handbook
September 2024

Contents

Introduction.....	3
What is AI?	4
Non-Generative AI.....	4
What’s New About Generative AI?.....	4
What Is an LLM?	4
Major LLMs.....	5
AI in the Publishing Industry.....	7
Existing Uses	7
New and Generative Uses.....	7
Audiobooks	7
Translation and Illustration	7
Marketing and Publicity.....	8
Within Agencies	8
Other Uses.....	8
AI and the AAA	9
Our Position on Generative AI.....	9
Our Goals.....	9
Who We Work With	9
AI in Politics and Legislation (last updated September 2024)	10
AI in the UK.....	10
AI in the European Union	10
AI globally	11
Notable lawsuits.....	11
Notable UK deals.....	11
Navigating AI in our work.....	12
How We Talk About AI.....	12
Licensing AI Rights / Collective Licensing	12
AI in Contracts / AI in the Publishing Process	13
AI in Submissions.....	13
Clients Using AI in Their Writing.....	13
AI Translations.....	14
AI in Audiobooks.....	14
AI in Covers and Illustrations	15
Final Note on the Handbook.....	16
Online resources.....	17

Introduction

AI has emerged, seemingly out of nowhere, as a major disruptor to our industry and many others. But the publishing industry has faced disruptors before – from Amazon emerging as a new retailer, to digital audiobooks emerging as an entirely new way to consume books, all the way back to the printing press itself – and each time we have adapted to the new technology, incorporated it into our workflows and learned to better serve our clients and, ultimately, our readers.

This time feels different to many of us. The rise of AI, particularly generative AI, has caught us all by surprise with the speed at which the technology has advanced and at which it has been taken up by all parts of the publishing community, from creators to publishers to consumers. It is, however, just another technological advancement to be turned from a behemoth of unfamiliar terms and uses into a tool that can be used to improve our working lives, should we wish to use it.

A major stumbling block on the path to conquering AI is understanding it. The technology has moved faster than we can keep up with, the law is struggling to adapt and update, and sometimes it feels like we need a degree in advanced computing just to understand what's going on.

The purpose of this handbook is to serve as an accessible introduction to AI: what it is, where it's come from, how it's being used, and how it's likely to change our working lives. It will not be an exhaustive document, but it is intended to be a living one, updated regularly as technologies change. It has been designed to help you navigate a rapidly changing landscape and feel empowered to discuss AI with your clients and your partners in the publishing world, and to make informed decisions.

What is AI?

Non-Generative AI

‘Artificial Intelligence’ is a kind of software that has been around since the 1950s. What separates all AI, both traditional and emerging technologies, from any other kind of computer programme is that AI can update itself and adapt to changing inputs.

At its simplest, AI is not terribly different from any other kind of software, and many of us have been using simple AI for years. Google Maps uses AI to update routes in real time, social media uses AI to tailor feeds, and Microsoft Word’s built-in spellchecker uses AI to check your writing against its constantly-updating database of grammar and spelling rules. The key thing to know about these technologies is that they mostly stick to using their own data. Netflix’s recommendation system, for example, uses AI to compare shows you’ve watched to similar customers in order to use their preferences to predict what you might like to watch next.

What’s taken the world by storm is generative AI, which is an emerging branch of technology.

What’s New About Generative AI?

Rather than crunching data to output patterns or predictions, a generative AI crunches data to output *more data of the kind it was trained on*. For example, a non-generative AI model can be trained on millions of photos of x-rays and told which do and don’t show tumours in order to ‘learn’ whether the next x-ray it ‘sees’ has signs of cancer. By contrast, generative AI is built from the ground up with the intention of creating output based on, and similar in format to, its input. Feed it a million photos of x-rays and it will be able to generate convincing looking x-rays to your specifications. The same is true of words – if you feed an LLM generative AI 100,000 books, it will be able to ‘write’ you another.

What Is an LLM?

Large Language Models (frequently abbreviated to LLMs) are generative AIs that have been trained on large amounts of written material for the purpose of generating new written material. They are frequently grouped with ‘text-to-image’ models, which are similar generative AIs that have been trained on images paired with descriptions of those images in order to generate images from text. You can give an LLM a written prompt (e.g. ‘Write me a poem about apples in the style of Ted Hughes’) and it will generate a body of text for you, and you can give a text-to-image model a written prompt (e.g. ‘Draw me a picture of apples in the style of Quentin Blake’) and it will generate an image for you. These new technologies are the ones that have evolved rapidly in recent years, much faster than international law has been able to keep up with, and are now expanding into many areas of the creative industries. There are emerging AIs that can generate videos, songs and many other kinds of falsely ‘creative’ outputs, all of which have been trained on huge amounts of real creators’ real work.

LLMs underpin most of these technologies – when you tell an AI to generate an image or a song, it’s an LLM that will process your written instructions and feed it through to the image

or music generation software. When you speak to an AI, AI software translates your spoken words into text and an LLM then processes that as though it had been written. Crucially, all LLMs must be trained on vast quantities of data – and by, ‘data’, we mean creative written work.

Books, articles, blog posts and social media comments are all examples of ‘data’ that can be used to train an LLM. At least one LLM is known to have been trained on the entirety of Wikipedia, among many other sources. Several major LLMs are known to have been trained on large quantities of books without the authors’ consent.

Many technologies utilise LLMs – it would be incredibly inefficient if everyone who wanted to add a chatbot to their website had to train their own from scratch. When a new piece of technology emerges, we encourage readers to ask which LLM that technology was based on, and what data was used to train it.

You may also encounter the term ‘model agnostic’ - this means that the software is happy to run on any LLM that the designers plug into it, rather than being tied into any one LLM.

Major LLMs

GPT

Generative Pre-trained Transformer, or ‘GPT’, is a series of LLMs created by OpenAI, an American AI research organisation. It has been through several major upgrades and at time of writing, the most advanced model is GPT-4. Its chatbot, ChatGPT, allows users to access and ‘talk’ to the GPT-4 LLM. GPT forms the basis for many other technologies and is the most widely available LLM.

BERT, PaLM, Bard and Gemini

These are all LLMs created by Google. BERT is one of the older LLMs and underpins a lot of Google’s later technology. PaLM (Path Language Model) is an older LLM that was designed to give more accurate and insightful responses. Bard is an LLM trained mostly on creative work for the purpose of generating work that feels like it has artistic flair. Gemini is Google’s current leading consumer-facing LLM and is notably multimodal – meaning its inputs and outputs can take the form of more than just text. It can process text, code, audio, images and video presented to it by the user and produce outputs in a similarly broad variety of formats.

Claude

Claude is a family of LLMs created by American startup Anthropic, first released in 2023. Anthropic’s goal is to produce an LLM that is “safe, reliable and useful to the public”, and the company was founded by seven former employees of OpenAI. Unlike other LLMs, including OpenAI’s GPT, Claude was trained using what Anthropic call ‘Constitutional AI’, a method designed to prioritise user safety.

Llama

This is the LLM designed by Meta (formerly Facebook) and is one of the few proven to have been trained on “Books3” – a leak of 196,640 pirated books by authors including Stephen King, Margaret Atwood, and Zadie Smith. Books3 was originally part of a much larger dataset called The Pile, which anyone wanting to train an LLM could use. Books3 was removed from The Pile after a legal complaint by anti-piracy group the Rights Alliance, but remains widely accessible online.

Grok

Developed by xAI – the ‘x’ here being the social network formerly known as twitter – Grok is an LLM distinguished by being exclusive to X. It is constantly updating itself in real-time from everything posted to X.

Other LLMs

While the LLMs listed above are the ones you’re most likely to encounter on a daily basis, many others exist, including Cohere (which was designed with business needs and instruction rather than chat in mind), Falcon (which was designed to feel like a natural chat with a human), and others such as DBRX, Phi-3 and Mixtral. New LLMs are being trained often, and older LLMs are being updated and occasionally renamed.

AI in the Publishing Industry

Existing Uses

It is worth reiterating that AI isn't new, and that it's been integrated into our workflows for a long time. Word's spelling and grammar function, for example, is an AI programme we all encounter on a daily basis. Excel uses AI in many of its functions, including the way it formats cell input and autocompletes formulae, and translation software such as Google Translate are built on AI. None of these programmes harm creatives or infringe copyright.

New and Generative Uses

There are many ways that AI technology can and is being incorporated into our working lives. Some of these emerging technologies use generative AI and some don't, some are designed to save time on administrative work and some are in direct competition with creatives. Several publishers are using AI tools internally to ease the admin burden on their staff. Several of these, including [Hachette UK](#) and [Pan Macmillan](#), have stated commitments to use AI responsibly while putting creatives and their outputs first. The nature and function of the technologies being trialled has not yet been made public.

Audiobooks

Generative AI is being used by several publishers to produce audiobooks without the expense of a human narrator and the studio time that audiobooks generally require. Text-to-speech technology can produce audiobooks in around an hour, in a multitude of languages and with a range of voices and styles available. Some publishers, such as [HarperCollins](#), are using the technology to produce audiobooks for backlist titles in non-English markets that would not receive audio versions otherwise. They argue this allows them to focus on human-narrated frontlist titles. Other companies, such as Urano World, are moving to an AI-first approach to their audio publishing. This technology can also be used by self-published authors to produce audiobooks of their own work.

Translation and Illustration

Generative AI is also being used to translate and illustrate works. [An April 2024 SoA survey](#) revealed that a third of translators and a quarter of illustrators reported having lost work to AI. It's not clear which publishers are using AI as a matter of policy in these cases, but there have been several high-profile cases of book covers being demonstrated to contain AI-generated assets. It is worth noting, however, that cover designers often use stock image libraries to source assets from which to build a cover, and it's not always immediately clear which images in those libraries are created by artists and which are generated by AI.

Marketing and Publicity

Generative and non-generative AI are being employed in marketing and publicity, with companies such as [Shimmr](#) using artificial intelligence to create marketing materials and to automate advertising. AI-powered advertising campaigns can continuously update themselves, using generative AI to create new copy and images, in response to the performance of their adverts being continuously monitored by non-generative AI.

Within Agencies

AI companies are creating tools for authors that aim to use emerging technology to critique and develop their work. Companies are also producing software aimed at helping agents to manage their unsolicited submissions, claiming to be able to ‘analyse’ manuscripts and produce reports on quality, content and genre. These technologies are in their infancy, and it remains to be seen to what extent they’ll be adopted. It is notable that while the technologies being used to review manuscripts may not use LLMs or generative AI, the responses that the software will produce based on those review almost certainly will.

Other Uses

AI is also being utilised within bookselling. German bookshop chain [Thalia](#) has partnered with an AI company to improve product discovery and optimise their book recommendations, allowing them to double the range of titles they offer while cutting down on the admin that would usually be required to make that profitable.

More broadly in the books industry, digital library [Perlego](#) has developed an AI tool aimed at optimising research, responding to queries using contextual understanding rather than a simple keyword search and suggesting more relevant books and excerpts than a traditional search would be capable of.

All of the examples listed are emerging technologies with the potential to optimise or disrupt different parts of the industry. It is notable that not all of these require the use of LLMs trained on copyrighted material, and that not all of them necessarily use generative AI.

AI and the AAA

Our Position on Generative AI

It is the opinion of the IP, Copyright and AI subcommittee that, as things stand and in the absence of a collective licensing agreement, there is no way to use generative AI, whether as an author or as a publisher, without infringing the rights of countless creators. We believe we cannot protect the work of our clients whilst also endorsing AI.

Our Goals

The purpose of the Sub-Committee on AI is to act as an educational resource for the AAA membership. We intend to do this by:

- Filtering, synthesizing and sharing information from allied organizations across the creative industries and other sources, including monitoring and summarising ongoing legal cases worldwide
- Facilitating an interactive forum in which the membership can approach us to discuss new challenges and opportunities as they arise
- Creating and maintaining (this!) document of best practice for the reference of member agencies

Who We Work With

As an entirely voluntary organisation, the AAA lack the time and resources to navigate industry-changing developments like AI on our own. Thankfully, we are part of an expert network of organisations. The following all provide helpful resources, which may be helpful for your own reference and/or for your clients (see section “Online Resources” at the end of this document for a list of links).

The Society of Authors (SoA)

The Creators’ Rights Alliance (CRA)

The Publishers’ Association (PA)

The British Copyright Council (BCC)

The AALA and the ALAA

The Authors’ Licensing and Collecting Society (ALCS)

Additionally, our interests are represented in parliament by the lobbying group the **Alliance for Intellectual Property** of which we are a member, who advocate for 23 organisations spanning “representatives of the audiovisual, music, toy and games, business software, sports rights, branded manufactured goods, publishing, retailing, image, art, and design sectors.”

AI in Politics and Legislation (last updated September 2024)

As a rapidly developing technology, the majority of progress in the field of AI is happening ahead of legislation. Policymakers in the UK and further afield are working to create the necessary laws, checks and balances to encompass this new technology. For the creative industries, and for publishing specifically, the major concern is that AI companies currently ‘train’ their models on unlicensed ‘data’ (i.e. writing), illegally using copyright-protected works scraped from the internet without authorisation or permission.

By its nature as an online technology, AI is difficult to legislate on a national level and requires a global outlook. There are many active lawsuits, as well as umpteen committees, reviews and inquiries running internationally. Below is a distilled overview of those likely to be most influential.

AI in the UK

The UK has not yet passed anything into law concerning AI. An ‘AI Bill’ was in the works under the Conservative Government but was not passed before the 2024 General Election. At the moment, we are reliant on pre-AI copyright legislation to protect authors’ rights, but this is unfit for purpose and is being ignored by both domestic and international AI developers.

In the July 2024 King’s Speech, the new Labour Government set out its intention to place safety and legality requirements on those developing the most powerful artificial intelligence models, but didn’t commit to blanket legislation concerning AI. In July 2024, the Government appointed Matt Clifford, tech entrepreneur and Chair of the Advanced Research Invention Agency, to develop an AI Plan. During this process, Clifford will engage with academic, industry and civil society experts to “set out how the Government can support the growth of the AI sector, enabling it to compete globally” and “include actions designed to boost the responsible adoption of AI across all parts of the economy”. We don’t yet know when the Plan will be published.

AI in the European Union

The EU’s AI Act came into force on the 1st August 2024, becoming the world’s first AI-specific legislation.

The Act seeks to “ensure that artificial intelligence developed in the EU is trustworthy and protects fundamental rights [and] is intended to create a harmonized internal market for artificial intelligence in the EU, promote the use of this technology, and create a favourable environment for innovation and investment”.

The most relevant provisions are that AI outputs/AI creators will have to:

1. Disclose that the content was generated by AI
2. Design models to prevent them from generating illegal content
3. Publish summaries of copyrighted data used for training

Notably, point 3 does not prohibit companies from using copyrighted data, it only commits them to admitting when they have. It remains to be seen how this law will be enforced, if at all.

AI globally

The USA's AI policies are likely to be the amongst the most impactful on how AI is legislated globally. In April 2024, Democrat Congressman for California, Adam Schiff, introduced the "Generative AI Copyright Disclosure Act" which would "require a notice to be submitted to the Register of Copyrights prior to the release of a new generative AI system with regard to all copyrighted works used in building or altering the training dataset for that system. The bill's requirements would also apply retroactively to previously released generative AI systems, and it was widely endorsed by creative organisations in the US, including the Writers' Guild.

Notable lawsuits

There have been a number of high-profile lawsuits launched in response to unlicensed usage of creators' work by AI training models. Amongst the most significant are:

Getty Images vs Stability AI – this lawsuit was launched in Delaware in 2023 before coming before the High Court in the UK in 2024. The High Court judge who approved the case for trial said Getty's claim has a "real prospect of success" in relation to London-based AI firm Stability AI's Stable Diffusion model's "image-to-image feature", which Getty claims allows users to make "essentially identical copies of copyright works, notably using 12 million unlicensed images from its databases." The result of the case is likely to have significant weight in future legislation.

The New York Times vs OpenAI/Microsoft - the NYT was, in its own words, "the first major American media organization to sue the companies, the creators of ChatGPT and other popular A.I. platforms, over copyright issues associated with its written works." The Times alleges infringement of its, and its writers', copyrights on a mass scale, and is claiming the infringing LLMs should be destroyed. As with Getty vs Stability, the eventual judgement is likely to be highly influential on future legislation. However, at time of writing, the suit appears likely to be a long one, and may not be resolved for several years.

Notable UK deals

Large academic publishers have been at the forefront of licensing academics' work to AI firms for training data. Taylor & Francis (which owns Routledge) struck a \$10m deal with Microsoft in July 2024 and was [reported](#) to be due to make \$75m overall from 'AI partnerships'. In August 2024, Wiley was [reported](#) in the *Bookseller* to have made deals worth \$44m with various AI companies. In both cases, it is notable that the authors of the individual works involved had not been informed, consulted, offered the opportunity to opt-out, or paid.

Navigating AI in our work

How We Talk About AI

It's tempting to slip into the language that the tech industry tells us is the correct way to talk about AI, language which reduces authors' work to pieces of exploitable data, rather than human-made things with intrinsic creative worth.

At the AAA, we avoid using language that undermines the creative input of our clients' work. We avoid referring to bodies of writing or illustrations as 'IP', 'data' or 'data sets', instead we talk about our clients' writing and books. Similarly, we try to avoid anthropomorphising the machines. LLMs are not 'taught' and they don't 'learn', they also can't 'write' or 'create'. They are, in fact, 'programmed' and 'optimized'. They can't 'think', 'decide' or 'understand' when responding to prompts, they can only 'process', 'evaluate' and 'generate outputs'.

These small changes in the way we discuss this emerging technology help to reaffirm our position: that creative work is not training data, it is the product of human ingenuity, and something worth celebrating and protecting.

Licensing AI Rights / Collective Licensing

At the forefront of the conversation around the use of creative work in training AI models is the question of how to remunerate authors. The illegal dataset 'Books3', for example, contained 190,000 stolen books – if this *had* been legally licensed data, how would a company go about clearing rights with 100,000+ individual authors? We can't turn back the clock on AI, so how do we ensure a fair future for the humans whose creative works are being alchemised into 'data'?

While old contracts are unlikely to reference AI in any way, the right to use a book for training an AI model is now something that should be actively considered when striking a deal.

One option is to assign AI rights to the publisher (much as you might assign anthology and quotation rights or translation rights).

Another option is retaining those rights and licensing them to companies directly, as a permission, with no intermediary.

A third option is collective licensing through a collecting society like the ALCS. The ALCS pitched their ability to handle AI licensing to the creative community in June 2024, saying "we have consistently secured remuneration for writers where new technologies result in their works being used without permission: in the 1970s it was for photocopying, in the 2000s scanning and online use and in the 2020s, it's AI."

AI in Contracts / AI in the Publishing Process

Similarly to broader national and international law, publishing contracts must catch up with the technological developments in AI to adequately express, legislate and – where necessary – protect both parties.

We know that blanket prohibition of AI is counterproductive; its non-generative uses are widespread, helpful, and have been practised for years (see section “What is AI?”), and in its generative capacity, there are also many ways in which AI can potentially aid the publication process – whether that’s generating metadata/keywords, iterating on copy, or digesting a large body of text into bullet points.

The idea of AI for ‘creative’ uses such as editing, copywriting, translating etc, has broadly been met with disapproval by authors and other industry creatives and, as such, many publishers have publicly stated they don’t intend to use it in some such capacities (see section “AI in the Publishing Industry”).

The SoA makes the following recommendations to its author members:

Ask your publisher to confirm that it will not make substantial use of AI for any purpose in connection with your work – such as proof-reading, editing (including authenticity reads and fact-checking), indexing, legal vetting, design and layout, or anything else without your consent. You may wish to forbid audiobook narration, translation, and cover design rendered by AI.

AI in Submissions

It is difficult to estimate how many writers are using AI in their work. We know that it has become endemic among students for essay writing, but in fiction particularly the breadth of use is harder to gauge. AI-generated text often *feels* AI-generated, and unsolicited submissions written using AI are still highly likely to be rejected based on the quality of their writing rather than explicitly because an agent smells AI. But the LLMs are being constantly updated and improved, and AI-detection becomes more difficult with each generation of LLM.

AI detection tools exist but are not infallible. Free online tools are available, but do raise questions as to whether agents have the right to run submitted work through third-party software. Subscriptions are also available for more reliable software, but third-party issues and GDPR would remain a concern. There have also been cases of illustrators finding their genuine creative work being automatically flagged as AI-generated on sites such as Instagram, because enough people are using generative AI to copy their style. This is, of course, another example of AI being unreliable.

Clients Using AI in Their Writing

There is little to no copyright protection for AI-generated text. Most countries have no legal protection at all for computer-generated works in any format. The UK is an exception, with an obscure law originally drafted to allow the video games industry to mass-produce text in the days before generative AI. The law has never been tested. There are some legal grey areas here

– for example, if a user generates output using an LLM and then substantially edits it, they can claim to own copyright in it, but this and other legal grey areas have also not been tested.

A writer does not own the copyright in text they have generated using AI and is not legally permitted to licence that copyright to a publisher. It is therefore of paramount importance that authors are honest about whether AI was used in their creative process, and to what extent.

[An SoA survey in April 2024](#) found that a fifth of surveyed fiction writers and around a quarter of non-fiction writers had used AI in their work. This use can span a variety of types of engagement. Some authors may use an AI chatbot to workshop plot ideas and character development before genuinely creating their work, and owning the copyright in it. Some authors may use generative AI to write sentences or paragraphs, which puts them in a legal grey zone as far as copyright is concerned. Some may rely solely on generative AI for their writing and not disclose the fact.

AI Translations

The use of AI to translate books has been one of the earliest (and indeed most feasible) applications of the improvements in AI technology, growing out of existing and widely-used free software like Google Translate.

The SoA have warned authors to be vigilant in relation to their work being translated by AIs rather than human translators, and have suggested that preventative language be inserted in contracts with their primary publisher when translation rights are sold, and likewise with individual foreign publishers.

The SoA have also raised the issue of impact on translators' income, with many already being asked to 'improve' translations made by AIs for significantly lower remuneration than they would receive for translating from scratch.

As with all AI-generated output, there is little to no copyright protection in AI-generated translations. If a translation of a work is not protected by copyright, then it's possible that such a version of the author's work could effectively enter the public domain.

AI in Audiobooks

AI narration for audiobooks is another area in which technology has come relatively far, and continues to improve rapidly. Like its text counterpart, audio AI comes with copyright pitfalls - audiobook voice actors and performers are having their voices replicated by AI systems without acknowledgement, authorisation, or compensation, and also losing work to competing AI 'narrators'.

As with translation, the SoA have advised their members to prohibit use of AI narrators in their audio contracts.

AI in Covers and Illustrations

The Association of Illustrators (AOI) put out a statement in February 2023 regarding their stance on AI. The AOI believe that illustrations may only be used as AI training data if the copyright owner has given explicit permission for such use. Illustrators usually own the copyright in their work, and therefore the unauthorised copying of images then used commercially within a dataset can be considered a copyright infringement.

Illustrators may also have concerns about the commercial application of AI-generated works undermining their own existing working models, and the AOI are closely observing the extent to which businesses and commissioners will utilise this technology. For further information on navigating the changing landscape of AI in illustration, the AOI have put together a page with regular updates on the state of play, which can be viewed on their website [here](#).

On the flipside of this, many illustrators are keen to engage with new technologies and may wish to explore the use of AI in their own work. The position of the AOI is that they do not wish to hinder technological advancements or limit the ability of creative communities to utilise and benefit from AI, but instead hope that these platforms will be used ethically, and maintain that there needs to be a legal framework in place that does not freely exploit the intellectual property of creators without permission or remuneration. For this reason, the AOI does not intend to promote AI-generated artwork and works produced using AI text-to-image generative software are not currently eligible to be entered into the World Illustration Awards (presented by the AOI).

Final Note on the Handbook

The information provided in this handbook is for general informational purposes only. While we at the AAA have made every effort to ensure the accuracy and completeness of this information, the AAA and the handbook's individual authors can't assume any responsibility for any errors or omissions, and we encourage members to seek professional legal advice rather than basing any business decisions on this information. The content of this handbook is not intended to be a substitute for professional advice.

This handbook is not an exhaustive or definitive document. All information in this handbook is provided 'as is' at time of writing, and while we took every effort to provide accurate information, we cannot guarantee the completeness, accuracy or timeliness. The law and landscape are both liable to change rapidly.

Online resources

[Written ministerial statement](#) on the UK government's AI Plan

[EU AI Bill](#)

[ALCS AI principles](#)

[Society of Authors statement on AI](#)

[CRA statement on AI](#)

[PA statement on AI](#)

[AI at Pan Mac](#)

[AI at Hachette UK](#)

[AI an Penguin UK](#)