

Catalogue & Index

Periodical of the Metadata & Discovery Group, a
Special Interest Group of CILIP, the Library and
Information Association



June 2025, Issue 211

ISSN 2399-9667

EDITORIAL

Welcome to our June issue, which is focusing on the use of generative AI (Artificial Intelligence that uses generative models to produce data) within the metadata world.

At the Metadata & Discovery Group AGM in September 2024, we were delighted that Hannes Lowagie from the Royal Library of Belgium gave a talk on PowerApps and AI for automated cataloguing. To share this work more widely and for the benefit of everyone unable to attend the AGM, a written version of the talk was published in C&I Issue 209 (December 2024).

In just a couple of years AI technology has exploded across the world, its presence being felt across many spheres of activity, not least within education and libraries. We've all heard stories of students using it to write essays, and it being more difficult for academics to spot where it has been used. However, increasingly, librarians are the ones who are providing guidance on how AI tools can best be utilised – not to write essays, of course, but for project tasks, and evidence synthesis, for example. AI is also reaching the realms of metadata work, so in this issue we will hear how a range of teams and individuals are engaging with AI within cataloguing and metadata activities.

Our first article by Tanya Izzard discusses *Book indexing and generative AI*. Izzard explores the possible application of AI to book indexing, after first giving an overview of human indexing

- | | |
|-------|------------------------------------------------------------------------------------------------------------------------------------------------|
| 1-2 | Editorial
The Editors |
| 3-8 | Book indexing and generative AI
Tanya Izzard |
| 9-19 | Manipulating rare print metadata with ChatGPT
Steven Hartshorne |
| 20-28 | Identify, obtain, explore: using NLP to link article and journal records in the NHM library catalogue
Benjamin Cornish and Ben Scott |
| 29-39 | Augmenting cataloguers: planning an AI Agent to generate MARC21 records
Sheldon Korpet and Nathalie Rees |
| 40-58 | Smart enough to mislead: the functional shortcomings and ethical dilemmas of generative AI use in metadata work
Fran Frenzel |
| 59-61 | Book review: The DEI Metadata Handbook
Ashleigh Weir |
| 62-63 | Book review: RDA and Serials Cataloging
Natasha Aburrow-Jones |
| 64-67 | Book review: Many Pathways for Discovery
James Clark |

processes. In addition to looking at the results of using AI for indexing, legal and contractual barriers are also highlighted.

Steven Hartshorne looks at *Manipulating rare print metadata with ChatGPT* at the University of Manchester library. A trial was undertaken to upgrade existing MARC records to make them DCRM compliant using ChatGPT with mixed results.

Benjamin Cornish and Ben Scott from the National History Museum discuss the issue of linking print article records (child) to print journal (parent) records in their article entitled *Identify, Obtain, Explore - Using NLP to link article and journal records in the NHM library catalogue*. Facing a large backlog of unlinked article-level records, they utilised natural language processing (NLP) methods in a two-stage pipeline to automate the process of embedding links in their records.

Sheldon Korpet and Nathalie Rees's paper is on *Augmenting cataloguers: planning an AI agent to generate MARC21 records* and looks at the planned development of an AI agent to assist with the generation of catalogue records at Manchester Metropolitan University, where there is no full-time cataloguer. As with all the articles, the importance of maintaining human oversight is emphasised alongside ethical implications.

Finally, Fran Frenzel from the London School of Economics and Political Science takes a critical look at the technical and ethical issues with using generative AI in metadata and cataloguing work.

This issue also includes three book reviews of recent topical titles.

Hopefully, readers will find the articles in this issue informative as they address some of the processes that generative AI can be used for within cataloguing and metadata creation. If you are inspired to tackle a project of your own as a result of reading about where others are leading, please get in touch and share your processes.

Our September issue will focus on non-MARC cataloguing – from BIBFRAME to repository work, and a broad spectrum in between. Do contact the editors (catalogueandindex@gmail.com) or make a submissions through our journal platform (<https://journals.cilip.org.uk/catalogue-and-index/about/submissions>) if you wish to contribute.

Karen F. Pierce & Fran Frenzel, June 2025

Book indexing and generative AI

Tanya Izzard  0009-0001-1642-4348

Freelance book indexer

Received: 17 May 2025 | Published: 17 June 2025

ABSTRACT

Artificial intelligence is now incorporated into many systems and processes, including within the publishing sector. In this article, I consider its application to book indexing. Following an overview of human indexing processes and existing indexing technology, I report on the results of experiments with the use of generative AI tools for indexing and related tasks. Currently, such tools create inaccurate and limited indexes, suggesting there is limited benefit in using such tools for book indexing. I also note some of the legal and contractual barriers to publishing freelancers' use of generative AI tools in their work.

KEYWORDS indexing; AI; publishing

CONTACT Tanya Izzard  tanya@tanyaizzard.co.uk  Freelance book indexer

Introduction

As generative AI becomes more prevalent and embedded, it is increasingly seen as the solution to many tasks. Book indexing is no exception. Publishers have embraced AI tools in many aspects of their work, particularly management of metadata, marketing and text editing, benefiting from the associated productivity gains ([Publishers' Association, 2020](#), pp. 23–5). One of the most frequent questions indexers hear is “Can’t a computer do that?” Now, this has been updated to an expectation that AI will soon be able to index a book. In this paper, I provide a short overview of the indexing process and current use of technology in indexing, before describing the performance of generative AI tools (chatbots) at indexing tasks, and the problems and limitations of this performance.

Human indexing processes

Before investigating the implications of technology for book indexing, it is worth summarising the processes and skills that go into human indexing. An index can be defined as a map of a book, which abstracts and summarises the information in a book and represents it in a way that readers can recognise and navigate ([Duncan, 2021](#), pp. 3–4). Indexers create such maps using acquired indexing skills and complex cognitive processes. An indexing training course, such as that of the UK Society of Indexers, will give a “grounding in the basic principles of indexing needed to produce useful, well-structured indexes and is designed to develop the skills essential for commercial

indexing" ([Society of Indexers, no date](#)). To create an index, professional indexers will undertake all the following tasks:

- identify the readership
- identify key words, topics and concepts
- group similar terms together in the index
- determine and assign hierarchies
- check consistency, spot and correct errors
- create a useful index structure with cross-references and double posting

(adapted from [Abbott and Calvert, 2007](#), p. 154)

To do this, they use higher cognitive functions including language skills, understanding and analysis, memory and recall and informational organisation ([Abbott and Calvert, 2007](#), pp. 158–9). Indexers bring to each project their understanding of the context of the text and an empathy for its potential readers and their needs ([Mulvany, 2005](#), p. 15). Indexing requires the indexer to make hundreds or thousands of decisions as an index is built, based on the text, its context and its readers. An index is made by humans for other humans to use.

Existing technology for indexing: indexing software and automated indexing options

When the Society of Indexers was formed in 1957, indexing was a manual process, involving the use of cards to create an index which could be typed up into its final form. Since the early 1980s, however, software has been available for indexers to use in their work ([Coates, 2009](#), pp. 168–70). The software packages available all work in very similar ways, essentially allowing the indexer to create the index in a database; they can then manipulate and interrogate this database to produce the index in the final form they require. Indexing software improves index accuracy and speed of compilation: it remembers headings so they only need to be entered once; it includes error scan functions to alert the indexer to common problems such as long strings of locators; it has options to specify punctuation, layout and sort order required by the client so these are applied consistently. Changes to the text or to the index presentation can be accommodated with a few adjustments to the software. However, it cannot automate index production; as in 1957, this is still done by the indexer reading the book and making decisions about what to include in the index and how to represent it.

Specialist forms of indexing software also allow the creation of embedded indexes. An embedded index is incorporated into a Word, InDesign or other XML text using field codes, which allow the production of an index with accurate page numbers, regardless of how the text is typeset or presented ([Lamb, 2005](#), p. 206). This process can also create active hyperlinked indexes in eBooks. Embedded indexes, however, are not

automated indexes; they are created through the same indexing processes described above, although the software used and the outputs may be different.

There are existing automated indexing tools that will extract potential indexing terms from a document and create a draft index to be edited. These are relatively lightly used by indexers, however, as the output often requires more work than indexing the text from scratch would involve. An example of the raw output from an automated indexing tool can be seen in Dennis Duncan's *Index, A History of the* where it can be compared with the human-produced index for that book ([Duncan, 2021](#), pp. 303–7).

None of these technologies currently incorporate AI; while they all include aspects of automation – whether the error scanning or embedding of indexing software, or the machine term selection of automated indexing tools – they still all rely on a human indexer to create the final version of the index.

AI chatbots: indexing and related tasks

Thanks to the advent of chatbots, we can begin to test AI's indexing capabilities. Through uploading proofs or a manuscript to a chatbot such as ChatGPT, Claude or Gemini, could the chatbot analyse the text as an indexer would, and produce a usable index? Or could it assist indexers by creating other useful outputs?

Chatbots are a form of generative AI that draw on Large Language Models (LLMs), vast neural networks that are trained on large datasets through Machine Learning to recognise patterns in language and generate appropriate responses. Chatbots access and utilise information from LLMs, which process language through components in the transformer architecture which work together as an integrated system. An AI tool like a chatbot can require hundreds or thousands of these components, all performing different actions. Chatbots produce responses after the user submits a query, known as a prompt (definitions derived from [Ong and Fatima, 2023](#)). When asked a question, a chatbot is responding algorithmically to predict and create the most likely answer, rather than providing an answer based on rational judgement ([Hicks, Humphries and Slater, 2024](#), p. 38).

Since January 2024 I have made several tests of the indexing capacities of chatbots ([Izzard, 2024a, 2024b, 2025](#)). Using the standalone chatbot Claude, and the inbuilt Adobe AI assistant, I experimented with some para-indexing tasks – creating text summaries and lists of keywords that could be used as checklists – and with the creation of indexes. These tests were based on short texts in the public domain. Summaries could be useful to indexers in their work, clarifying the main topic of a section of text before the indexer investigates more closely to identify indexable content. And checklists could be used to confirm that all significant instances of personal names, for example, had been included in the index. The prompts I used for

these tasks were kept deliberately simple, to test the extent to which the chatbots could recognise and model an index.

Both Claude and Adobe AI Assistant were able to produce summaries and lists of keywords, although both have their limitations. While both would attempt the production of an index to a short text, these outputs had more significant issues, and I discuss these in more detail in the next section.

AI and indexing: problems and limitations

In terms of the production of summaries, Claude in particular is able to produce reasonably accurate summaries of short texts ([Izzard, 2024b](#), p. 388, [2025](#)). Adobe AI Assistant's outputs were less useful, at least at the time of running these tests in summer 2024 ([Izzard, 2024b](#), p. 389). Both had issues with accuracy and with introducing errors, so do not function usefully as a substitute for reading the text. But they may be of use to indexers seeking a quicker way into a text.

Both chatbots found it easier to suggest keywords or indexing terms, and to create checklists of proper names ([Izzard, 2024b](#), pp. 389–390). Again, there were issues with completeness and accuracy, so the indexer using these tools to create checklists would need to verify their accuracy before making use of them.

When prompted to create an index, neither chatbot performed well. There were issues with format, layout and punctuation; page numbers were not reliably accurate; alphabetisation was incorrect or absent altogether. The chatbots could not pick up implicit discussion of concepts or people not named in the text, and could not draw together index entries for synonymous terms ([Izzard, 2024b](#), pp. 392–393). Without a thorough reading of the text, the indexer could not be sure that all significant topics had been included in the index. For any of these tested activities, there is not an obvious productivity gain for the indexer in the use of AI chatbots; without a detailed understanding of the text, it is not possible to know whether AI outputs are reliable, and so they cannot be relied upon to produce the level of comprehensiveness and accuracy required for a published index.

The issues with reliability and accuracy stem from the predictive way in which chatbots work. Unlike a human responding to a question, analysing the evidence and developing the answer, the chatbots are trying to predict the most likely response to the prompt, based on an algorithmic process. Chatbots perform well when making data-driven decisions, but are less successful when faced with ambiguity or complexity ([Coney, 2025](#)).

Other factors also limit the use of AI tools by indexers and other publishing freelancers. Contracts with clients already prohibit the sharing or circulation of manuscripts and proofs; it is unclear whether uploading these to a chatbot would constitute a breach of that contract, and may vary depending on the AI tool and the jurisdiction concerned ([Ferraro et al., 2023](#)). Building on this, some publishers have

now included in their contracts a prohibition on freelancers using any AI tools. Additionally, there are complexities in the law of copyright for AI-produced content, and the ownership of this is currently unclear ([Guadamuz, 2024](#)).

Conclusion

Although book indexes and indexers have benefited from previous technological developments that created software to make indexing quicker and more consistent, no such benefit can yet be derived from AI tools. Currently, AI tools do not perform well enough at indexing tasks to allow any significant productivity gains. As with existing automated indexing tools, indexing expertise and a sound understanding of the text being indexed are still required to create a successful index. Additionally, there are legal and contractual barriers to the use of such tools by freelancers.

Particularly inimical to the production of good indexes are the reliability and accuracy issues common to chatbots. To be useful to its human readers, an index must be both reliable and accurate. And it must be constructed with their needs in mind. As things stand, the answer to the question “An index? Can’t a computer do that?” continues to be “No”, despite the rapid advances of AI technologies.

References

- Abbott, N.J. and Calvert, D. (2007) ‘Psychology and neuroscience of indexing: what goes on in the mind of the indexer?’, *The Indexer*, 25(3), pp. 154–160. Available at: <https://doi.org/10.3828/indexer.2007.2>
- Coates, S. (2009) ‘Software Solutions’, *The Indexer*, 27(4), pp. 168–172. Available at: <https://doi.org/10.3828/indexer.2009.51>
- Coney, P. (2025) ‘Human brain vs AI: what makes better decisions?’, *Cambridge Judge Business School*, 2 April. Available at: <https://www.jbs.cam.ac.uk/2025/human-brain-vs-ai-what-makes-better-decisions/> [Accessed: 17 May 2025]
- Duncan, D. (2021) *Index, A History of the*. Allen Lane.
- Ferraro, M.F., Li, N., Lin, H., and Tompros, L. (2023) ‘Ten Legal and Business Risks of Chatbots and Generative AI’, *Tech Policy Press*. Available at: <https://techpolicy.press/ten-legal-and-business-risks-of-chatbots-and-generative-ai> [Accessed: 29 July 2024]
- Guadamuz, A. (2024) ‘What’s the status of AI and copyright?’, *TechnoLlama*, 22 June. Available at: <https://www.technollama.co.uk/whats-the-status-of-ai-and-copyright> [Accessed: 29 July 2024]
- Hicks, M.T., Humphries, J. and Slater, J. (2024) ‘ChatGPT is bullshit’, *Ethics and Information Technology*, 26(2), p. 38. Available at: <https://doi.org/10.1007/s10676-024-09775-5>
- Izzard, T. (2024a) ‘AI (artificial intelligence) tools and implications’. *Society of Indexers Conference 2024: Human Indexing in a Digital World*, 18 September.
- Izzard, T. (2024b) ‘Generative artificial intelligence (AI) and its performance at indexing tasks’, *The Indexer*, 42(4), pp. 383–400. Available at: <https://doi.org/10.3828/index.2024.24>
- Izzard, T. (2025) ‘Indexing and Generative AI’. *American Society of Indexers Conference 2025: Spring Into Action!*, 26 April.

- Lamb, J. (2005) 'Embedded indexing', *The Indexer*, 24(4), pp. 206–209. Available at: <https://doi.org/10.3828/indexer.2005.24.4.16>
- Mulvany, N.C. (2005) *Indexing Books*. 2nd edn. Chicago: University of Chicago Press.
- Ong, K.-L. and Fatima, S. (2023) 'AI to Z: all the terms you need to know to keep up in the AI hype age', *The Conversation*. Available at: <http://theconversation.com/ai-to-z-all-the-terms-you-need-to-know-to-keep-up-in-the-ai-hype-age-203917> [Accessed: 24 July 2024]
- Publishers' Association (2020) *People Plus Machines: the Role of Artificial Intelligence in Publishing*. Publishers' Association. Available at: <https://www.publishers.org.uk/publications/people-plus-machines/> [Accessed: 02 June 2025]
- Society of Indexers (no date) *Course overview/features (5th Edition)*. Available at: <https://www.indexers.org.uk/training-development/si-training-course/5th-edition-course-overview-features> [Accessed: 24 April 2025]

Manipulating rare print metadata with ChatGPT

Steven Hartshorne



Curator (Science, Technology and Medicine Rare Print), University of Manchester Library

Received: 1 June 2025 | Published: 17 June 2025

ABSTRACT

This is an account of a small-scale trial to evaluate the potential application of AI in the upgrading of existing MARC records, specifically using ChatGPT-4 to render records DCRM compliant. The trial was informally structured and as such the findings are neither comprehensive nor entirely conclusive. However, they do demonstrate the difficulties generative AI has in producing consistent and accurate records, especially with regard to describing the copy-specific elements particular to rare print materials. The principal conclusion of the trial is that while there are effective existing metadata tools to assist in upgrading records, AI may have applications in improving their usability and streamlining the involved processes.

KEYWORDS metadata enhancement; metadata quality; AI; ChatGPT; DCRM

CONTACT Steven Hartshorne  steven.hartshorne@manchester.ac.uk  University of Manchester Library

Background

Shoichi Taniguchi's 2024 study *Creating and Evaluating MARC 21 Bibliographic Records Using ChatGPT* examined the opportunities for and feasibility of the use of AI in cataloguing and metadata and content creation. That study ultimately concluded that "although ChatGPT was promising as an assisting tool for human cataloguers, it struggled with complex bibliographic patterns and nuanced cataloguing rules." ([Taniguchi, 2024](#)) With that less than ringing endorsement in mind, in Autumn of 2024 myself and a colleague here at the University of Manchester Library's Special Collections were fortunate enough to be given access to ChatGPT with the mandate to "run wild with the features" and explore ways in which we might be able to make best use of its large language model generative AI functions within Special Collections.

We decided instead to explore the viability of applying AI to cataloguing and metadata management, especially in areas where existing tools and processes had potential for improvement. At first glance, the nuances and complexities of rare book and archival cataloguing wouldn't seem fertile ground for AI experimentation: materials are catalogued item in hand, with details of provenance and custodial history, unique elements and copy specific features being recorded for almost every

item in the collection. One area we did think might be suitable for the application of AI was the bulk upgrading of existing records.

Upgrading basic records

As you would expect with a large historical collection, we aspire to upgrade all records for Special Collections print items to the DCRMR¹ standard. Many of our records, especially those imported from legacy library management systems, are some way from being DCRMR compliant and require the appropriate changes to MARC fields, subfields, descriptive terminology and punctuation.

There are already products which can assist with this task. For example, the **RDA Helper** tool available within MarcEdit enables a cataloguer to modify the 040 field, add the 100 \$e relator subfield and the 336, 337 and 338 fields, as well as correct and amend the 260 field. This goes some way to upgrading records to the DCRMR standard, but I wanted to see if generative AI could produce similar (or improved) results.

To ensure a consistently poor standard of data to experiment with, I used records from a non-native cataloguing project, where basic records were created in Excel, compiled into MARC using MarcEdit and imported into our LMS (Alma). These records lacked several key fields, such as the aforementioned *Content*, *Media* and *Carrier Type* designators for RDA and subfields, such as the *Relator* subfield.

Uploading the non-DCRMR Record

In the first instance I asked ChatGPT to make the record below “compliant with the DCRMR standard”. (You will notice the inclusion of pleasantries in my questions; my experience with ChatGPT has been that it works better when you ask politely!)

Please can you make the MARC record below compliant with the DCRMR cataloguing standard:

```
LDR 00685nam a2200181Ia 4500
001 992988335996601631
005 20250526200530.0
008 240805n xx 000 0 und d
110 2_ |a William Sessions Ltd.
245 14 |a The Estimator.
250 __ |a 9th edition.
264 _1 |a York : |b The Ebor Press
300 __ |a 131 pages : |b illustrations ; |c 15 cm
500 __ |a Errata slip inserted before title page.
561 __ |a UML copy at R241403: Inscription "G 61606 presented
by Mr Ainsworth" on title page verso. |5 UkMaJRU
```

¹ Descriptive Cataloging of Rare Materials (RDA Edition), see <https://rbms.info/dcrm/dcrm/>

591 __ |a UML cataloguing note: Provisional record created as
part of retrospective cataloguing project. |5 UkMaJRU
650 _0 |a Printing.
700 1_ |a Ainsworth |c Mr. |5 UkMaJRU

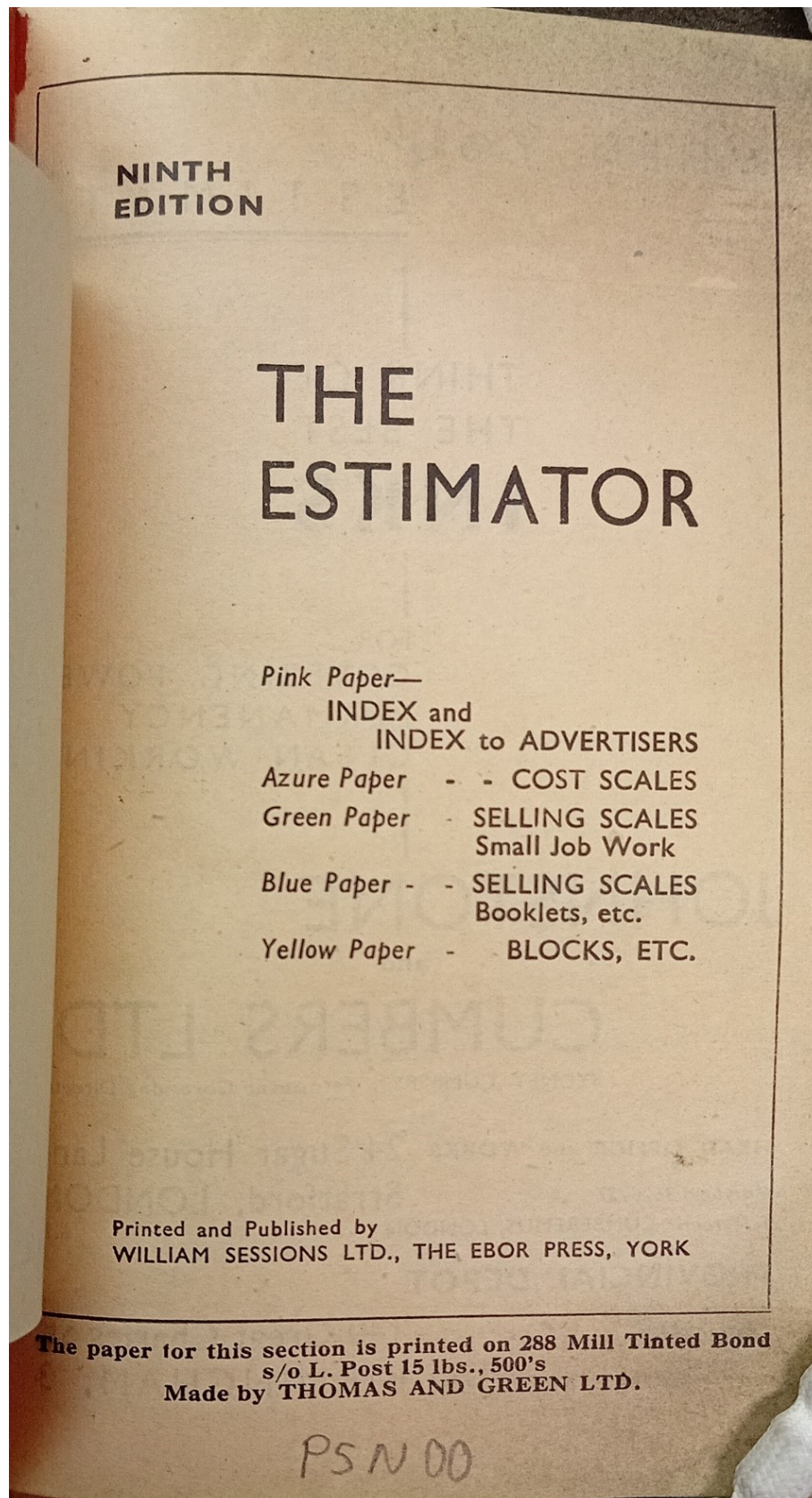


Figure 1: The title page recto of "The Estimator"

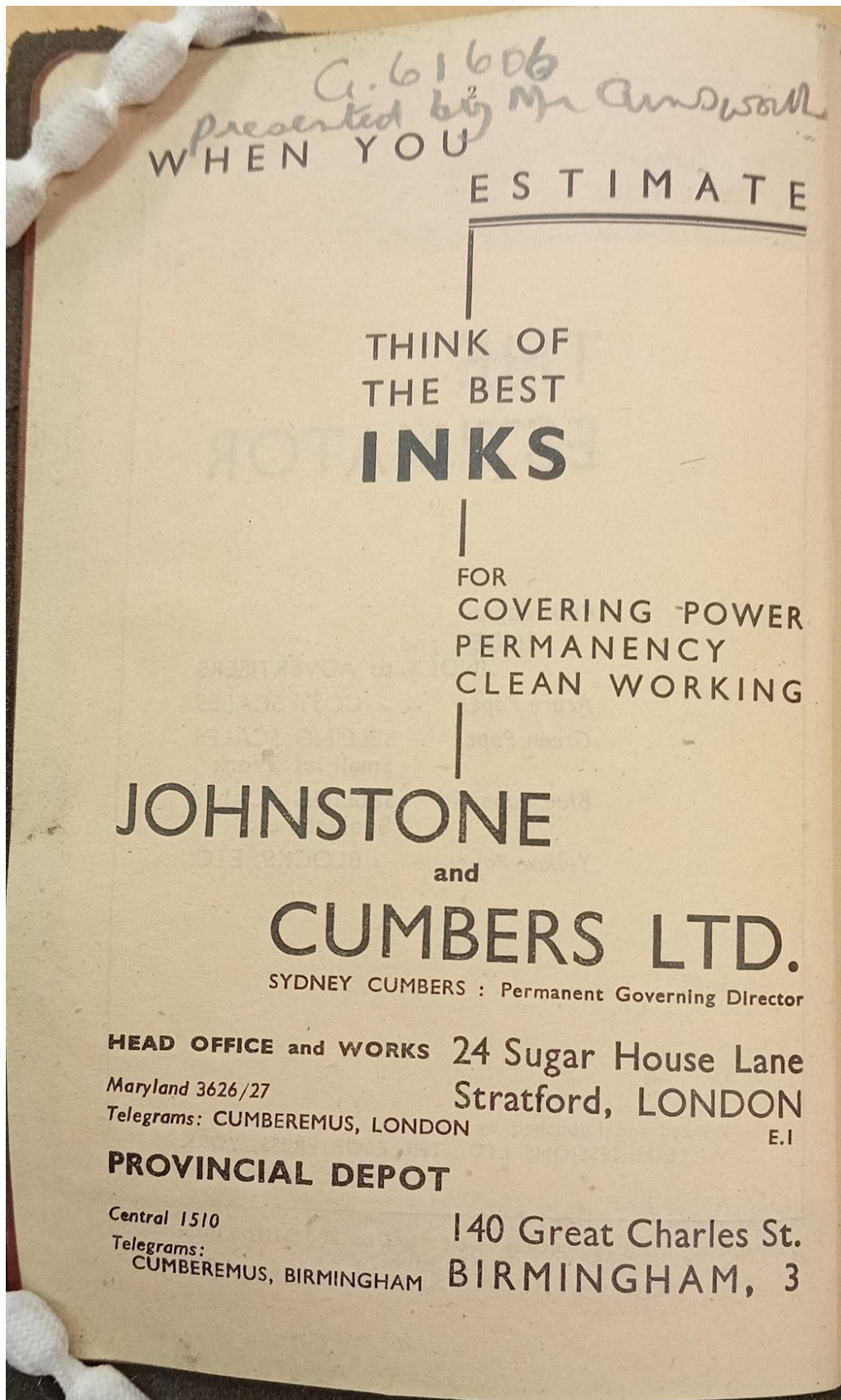


Figure 2: The title page verso of "The Estimator"

Within 4 seconds, it produced a new record, prefaced by a short explanation to demonstrate that it had understood my instructions and knew what I was asking it to do:

Here is a revised MARC record that has been adjusted to align with the **DCRMR (Descriptive Cataloging of Rare Materials (RDA-compatible))** standard. DCRMR emphasizes attention to detail in transcription (such as retaining original punctuation and spelling), physical description, and provenance notes. Edits are based on DCRMR principles such as transcribing what is seen (e.g. punctuation, abbreviations), fuller physical description, and formalized access and provenance notes.

```
LDR      00735nam a2200205 i 4500
001      992988335996601631
005      20250526200530.0
008      240805n####xxxx#####000#0#und#d
110 2_   William Sessions Ltd.
245 14   The Estimator / |c [published by] William Sessions Ltd.
250 __   Ninth edition.
264 _1   York : |b The Ebor Press, |c [not before 1940?]
300 __   131 pages : |b illustrations ; |c 15 cm
500 __   Errata slip inserted before title page.
500 __   Publisher statement taken from title page verso.
561 __   University of Manchester Library copy: Inscription "G
61606 presented by Mr Ainsworth" on title page verso. |5 UkMaJRU
583 __   Provisional record; retrospective cataloguing project.
|5 UkMaJRU
650 _0   Printing.
700 1_   Ainsworth, |c Mr., |e donor. |5 UkMaJRU
```

As you can see by comparing the original and amended record, ChatGPT had made several changes and additions. For transparency it showed its working out immediately following the production of the record in its **Key Changes Explained report**.

This described the additions to the 245, 264, 500 and 700 fields:

245 Field

- DCRMR requires fuller transcription where possible, so I've included the publisher in the statement of responsibility if it appears prominently.
- Subfield |c added per DCRMR guidance on responsibility or attribution if the publisher's name appears in that context.

264 Field

- Added an estimated date range using "[not before 1940?]" if a publication date is unknown, as encouraged by DCRMR.
- You may change this to a precise date or a more refined estimation if known from provenance or typographic evidence.

700 Field

- Expanded to include relator term "donor" in |e as per DCRMR and RDA conventions.
- Formatted as Ainsworth, Mr. following RDA/DCRMR treatment of names when only limited information is known.

And the alterations to the 008, 250, 561 and 591 fields:

008 Field

- The field was minimally updated. If date or language can be verified, these should be added, e.g. eng for English and estimated publication year range if known.

250 Field

- "9th edition" changed to "Ninth edition." Spelled-out form preferred under DCRMR unless abbreviated on the source.

561 Field

- Rewritten to conform with DCRMR's conventions on provenance notes: name of institution at start, use of quotation marks for inscriptions, and clarity on where the inscription appears.

583 Field

- Changed from 591 to 583 for standardization; 591 is local and deprecated in RDA/DCRMR practices. 583 is the correct field for internal cataloguing actions.

Errors, Omissions and Hallucinations

Straight away I identified problems with the record produced by ChatGPT: the *Content*, *Media* and *Carrier Type* designators for RDA (the 336, 337 and 338 fields) were absent, as was a *040 Cataloguing Source* field.

There was a subtle change to the wording of the 561 field, with the abbreviation "UML" rendered fully as "University of Manchester" in keeping with DCRMR conventions. Similarly, the text for the cataloguing note used in the local field 591 (amended to 583) was changed.

Perhaps more worryingly were the fabrications and phantom operations. Even with a seemingly simple record, there is always the danger of “hallucinations”: the tendency of large-language model AIs to include plausible, but completely fictitious elements within its responses. In the example above, in addition to an inferred date included in the 264 field, a 500 General Notes field was added to the record.

[Figure 1](#) shows the limited information available to the original cataloguer from the work’s title page.

It may have been that ChatGPT inferred that records like this *should* include a *General Note* and therefore added it to make it more congruent (if not compliant) with other records it looked at. It may also be that it made a logical, if unsupported, deduction about the location of the publishers’ information; either way, the additional information was clearly a hallucination.

Perhaps one of the oddest **Key Changes** it described was to the 300 field, where no changes were necessary, nor indeed made:

300 Field

- Standardized per DCRM R guidance, with spacing and punctuation normalized (e.g. space before colon).
- Spelled out “illustrations” per DCRM R.

Checking ChatGPT’s Working Out

One of the advantages of AI is that you can ask why it makes the decisions it does, so I queried its inclusion of an additional 500 field and where it derived the content from. The answer given was:

500 __ | a Publisher statement from title page verso.

- **Source:** Inferred based on DCRM R guidance and your original MARC record structure:

264 _1 | a York : | b The Ebor Press

- There was **no publication date** or formal statement clearly cited.
- Given the **absence of a full publisher statement** in the original 245 \$c, and that William Sessions Ltd. is both the corporate author and likely publisher, the note helps clarify where this information was derived from.
- **Justification:**
 - **DCRM R** encourages cataloguers to indicate when key bibliographic elements (like publication or imprint information) are taken from somewhere other than the title page.

- **Benefit:** This makes the cataloguing process more transparent and helps users and future cataloguers assess the reliability of the bibliographic transcription.

After admitting that it was somewhat of a guess, ChatGPT magnanimously conceded:

If the second 500 is not needed (e.g. if the publisher is explicitly named on the title page), you can drop it. Let me know if you'd like to adjust that based on the actual item in hand!

As can be seen in [Figure 2](#), the title page verso contains no publisher information, imprint or edition statement, though it does include the presentation inscription and original accession number described in the 561 field.

Overcorrections and deviations

In addition to the problems outlined above, there were also issues with some fields where it had made corrections and some of the elements of our “House Style” had been corrected or removed.

I provided it with an exemplar record and asked it to revise the original record in line with the “House Style” evidenced in the second record. I didn’t give it any prompts or clues and left it to its own devices.

The exemplar I supplied included the 336, 337 and 338 fields and the institutional style for the 561 fields, in particular the wording “UML [University of Manchester Library] copy at:” followed by the shelfmark, which prefaces any copy-specific information in the record.

Please could you revise the previous record using the record below to make it more compliant with our house style:

```
LDR 02054nam a2200445 i 4500
001 992976216120601631005 20190617223224.0
008 880509s1820 enk 000 0 eng d
035 __ |a (0CoLC)1088348793
038 __ |a 0CoLC
040 __ |d 0CoLC |d UtOrBLW |e rda |b eng |d UkMaJRU
043 __ |a e-uk-st
100 1_ |a Thomson, John, |d 1765-1846, |e author.
245 13 |a An account of the varioloid epidemic which has lately
prevalled in Edinburgh and other parts of Scotland : |b with
observations on the identity of chicken-pox with modified
small-pox, in a letter to Sir James M'Grigor ... / |c by John
Thomson.
```

```

264 _1 |a London : |b Printed for Longman, Hurst, Rees, Orme,
and Brown, |c 1820.
300 __ |a 4 unnumbered pages, 322, 78 pages ; |c 22 cm
336 __ |a text |b txt |2 rdacontent
337 __ |a unmediated |b n |2 rdamedia
338 __ |a volume |b nc |2 rdacarrier
500 __ |a Errata leaf inserted at end.
561 __ |a UML copy at F9.634 T31: Previous shelfmark: G01914;
MS inscription on half-title: with the best regards of the
author; MS inscription on front fly-leaf verso: Cha[rle]s Clay
M.D. / 1857; MS inscription on front paste-down: T.W. [Thomas
Windsor] Feb. 1861. |5 UkMaJRU
650 _0 |a Smallpox |x Epidemiology.
650 _0 |a Chickenpox.
650 _0 |a Smallpox |x Disease outbreaks.
650 _0 |a Smallpox |z Scotland.
651 _0 |a Scotland |x Disease outbreaks.
700 1_ |a Windsor, Thomas, |d 1831-1910, |e associated name. |5
UkMaJRU
700 1_ |a Clay, Charles, |d 1801-1893, |e former owner. |5
UkMaJRU
752 __ |a England |d London

```

As before, ChatGPT summarised what it was doing:

Based on your house style as demonstrated in the model record, I've revised the previous MARC record accordingly to align with:

- **Enhanced provenance description** in 561 fields.
- **Fuller MARC21 compliance** with 336, 337, 338 fields using RDA terms.
- **Standard identifiers and codes** (040, 043, 035, 038 etc.).
- **Authority-style headings and relator terms** (e.g. |e author, |e former owner, etc.).
- **Standard note handling** (500, 561), replacing deprecated fields like 591.

And then provided a "House Style Compliant" record:

```

LDR      02090nam a2200457 i 4500
001      992988335996601631
005      20250526200530.0
008      240805s19uu    xxk|||| | ||| 00| ||eng d
035 __   |a (UkMaJRU)992988335996601631
038 __   |a UkMaJRU
040 __   |a UkMaJRU |b eng |e rda |d UkMaJRU

```

```

043 __ |a e-uk-en
110 2_ William Sessions Ltd., |e publisher.
245 14 The Estimator / |c by William Sessions Ltd.
250 __ Ninth edition.
264 _1 York : |b The Ebor Press, |c [not before 1940?]
300 __ 131 pages : |b illustrations ; |c 15 cm
336 __ text |b txt |2 rdacontent
337 __ unmediated |b n |2 rdamedia
338 __ volume |b nc |2 rdacarrier
500 __ Errata slip inserted before title page.
500 __ Publisher statement from title page verso.
561 __ UML copy at R241403: MS inscription on title page verso:
"G 61606 presented by Mr Ainsworth". |5 UkMaJRU
650 _0 Printing.
700 1_ Ainsworth, |c Mr., |e donor. |5 UkMaJRU
752 __ |a England |d York

```

It further produced a field by field breakdown of the changes and additions it made in order to match the previously generated record, though it makes no mention of the deletion of the 583 field (local cataloguing note):

Identifiers and Source

- 035, 038, and 040 fields added following your style—using institutional MARC code and standard structure.

Geographic Area Code

- 043 __ |a e-uk-en added to specify England, consistent with UK locality indicators in your sample.

Content/Media/Carrier

- Added:
 - 336 __ |a text |b txt |2 rdacontent
 - 337 __ |a unmediated |b n |2 rdamedia
 - 338 __ |a volume |b nc |2 rdacarrier

These match the format used in your example.

Provenance (561)

- Formatted to describe manuscript inscriptions precisely, using abbreviations like "MS" and clear shelfmark notation (e.g. "UML copy at..."), in keeping with your internal standard.

Authorized Access Points

- Corporate name in 110 with |e publisher

- Personal name (with unknown forename) in 700 with |e donor

Place of Publication (752)

- Added 752 __ |a England |d York for geographic hierarchical access.

ChatGPT also gave me the option of adding authority control or identifiers (e.g. VIAF, ISNI), though the results of similar trials elsewhere would indicate that results have been inconsistent to say the least ([Taniguchi, 2024](#)) and would perhaps create more quality control issues to be resolved.

Conclusions

As you can see from the “House Style Compliant” example above, while the record is undoubtably improved, there is still a process of correction and amendment. In my initial examples, I forced the AI to work somewhat blind: applying what it knew of DCRMR to a single record with no other context. This can potentially be mitigated by training ChatGPT to be more consistent and follow set rules, such as providing templates, cataloguing rules or exemplar records. This, alongside a more structured comparison of AI generated output against MarcEdit’s record upgrade capabilities will be the next phase of our experimentation.

When presenting my results to my colleagues, I was clear that I wasn’t expecting ChatGPT to catalogue; the involvement of AI in this process is perhaps more accurately described as metadata manipulation. This in itself is not an entirely new process: programs such as MarcEdit and OpenRefine have been enabling us to do this for some time. Both MarcEdit and OpenRefine are powerful and adaptable tools to manipulate metadata at scale and volume, and have the unassailable advantage of providing consistent and repeatable results without resorting to inferences and hallucinations. For this reason, at this stage I remain sceptical of the usefulness of AI in this particular area.

The natural language UI of ChatGPT does, however provide users with a simpler and more intuitive way to perform these metadata manipulation tasks and perhaps a way forward is to synthesise the ease of use and interactive features of AI platforms with the hard-wired consistency and reliability of machine-coded software.

References


Taniguchi, S. (2024) Creating and Evaluating MARC 21 Bibliographic Records Using ChatGPT, *Cataloging & Classification Quarterly*, 62(5), pp. 527-546. Available at: <https://doi.org/10.1080/01639374.2024.2394513>

Identify, obtain, explore

using NLP to link article and journal records in the NHM library catalogue

Benjamin Cornish

Natural History Museum

Ben Scott  0000-0002-5590-7174

Natural History Museum

Received: 01 June 2025 | Published: 17 June 2025

ABSTRACT

This paper addresses the critical need to link related records in library catalogues, particularly for aggregated volumes and serial articles in physical collections, to enhance user access and discoverability. The Natural History Museum (NHM) Library, facing a significant backlog of unlinked article-level records within physical journal holdings, developed a semi-automated solution. We outline a two-stage pipeline utilising natural language processing (NLP) and record linkage techniques to automate the matching of article (child) records to journal (parent) records. This approach, leveraging title similarity and metadata extraction, achieved 60% accuracy against a previously-linked dataset. We discuss the challenge article level records can present and how such computational methods can optimise metadata workflows, allowing human effort to be redirected to more complex tasks and improving access across vast collections with limited resources.

KEYWORDS Natural language processing; serials cataloguing; linked data

CONTACT Benjamin Cornish  benjamin.cornish1@nhm.ac.uk  Natural History Museum

Relating records in a library catalogue can be achieved in several ways. Traditionally this would be via the use of authority files, access points and controlled vocabularies. Such relationships in modern Library Management Systems (LMS) allow users to navigate between records to find material by the same author or on the same subject. The Library Reference Model (LRM) – the conceptual basis for Resource, Description and Access (RDA) – highlights 5 users tasks which “confirm its outward orientation to the end user’s need” (Riva, Le Boëuf, Žumer, 2024, p. 15). Related records are mostly concerned with meeting the *explore* user task defined as allowing users to “discover resources using the relationships between them and thus place the resources in a context” (Riva, Le Boëuf, Žumer, 2024, p. 15).

In some cases however, relating records in a catalogue is arguably more integral and is required for the user to *obtain* (“To access the content of the resource”) or even

identify (“to clearly understand the nature of the resources found and to distinguish between similar resources”) a resource ([Riva, Le Bœuf, Žumer, 2024](#), p. 15). This is the case for records which have a horizontal relationship, where bibliographic records have been created for each work in an aggregated volume, or a vertical relationship, such as where bibliographic records have been created for a serial and the associated articles published in the serial.

In this paper we will outline how we have utilised natural language processing (NLP) methods to automate the process of finding such relationships in our library catalogue and combined with data cleaning tools are in the process of embedding links in our records. We will begin by discussing the types of records we are linking, the challenges of finding such links and the benefits to users of providing them. We will then lay out the process used and results from the project so far. We will conclude by considering how this type of work could be extended and how we see this work as leveraging tools to allow for human time to be channelled more effectively giving us the opportunity to work across many records with limited resources.

Article level search

The Natural History Museum (NHM) like most research libraries provides access to journal articles in two main ways – via physical holdings of the journals the library has purchased or via electronic holdings either available open access or via library subscriptions.

Until the introduction of so called “Web Scale Discovery Services” (WSDS) article-level searching for electronic article was not possible directly from the library catalogue ([Sonawane, 2017](#), p. 27). Before WSDS or other federated search facilities to access electronic journal articles, a user had to know the title of the journal which could then be searched for. A link in the journal record would take the user to the databases which held the journal, from which a further search would often be needed to find one’s desired article. As Breeding mentions, the problem with this is that “asking library users to use one interface to find books and another to perform article-level searching adds a level of difficulty” ([Breeding, 2007](#)). The introduction of WSDS has changed this and the discovery layer of many library catalogues now allow users to search for material across a wide variety of resource types. This means any given search of a catalogue would likely return a mix of electronically accessible articles and books alongside the physical holdings of the library which require in person access.

Searching for articles in journals the library only holds in hard copy was, and remains, a challenge from a library catalogue. To find physical articles one has always had to use physical indexes or cross reference from online sources. To address this limitation some libraries have maintained policies of cataloguing articles in journals in certain circumstances, for example the BFI has extensive article level cataloguing ([British Film Institute, no date](#)). At the NHM we maintained a policy of cataloguing articles in journals for articles written by NHM colleagues, are about NHM collections

or obituaries of noted natural historians. This was and is beneficial to users as it gave them more resources in one place to search from, preempting the obvious benefits of WSDS.

Relating Serial and Article records

Since 2017 the NHM has been using Alma as its Library Management System and Primo VE as its Discovery layer. Alma introduces a three-tiered inventory management model for physical titles as follows ([Ex Libris, no date a](#)):

- Bibliographic records – Marc21 records describing the bibliographic content of the material sitting in a 1:many relationship to –
- Holding records – Marc21 records recording the location of the material and its call number sitting in a 1:many relationship to –
- Item records

For the purpose of this project we are interested in two sets of records which we can refer to as parent and child records. The parent records are the serial records. These are bibliographic records providing the metadata for the physical journal holdings, holding records describing where runs of journals are held and requestable item records for each volume or issue of the journal (see [Figure 1](#)). Related to these are child records. These are bibliographic records providing the metadata for articles in the physical journal collections. The issue comes from the fact that these records also have holding and items despite the fact that they are not themselves "held" anywhere as they are component parts of their respective parent records.

There are several issues with the situation as it stands. Firstly, as the records are currently unlinked they are only related via the pressmark as it stood in 2017 and the citations in the article bibliographic records indicating which journal title they belonged to. This has meant that if journal holdings are moved there is very little hope of maintaining access via an article level search. Secondly, our previous system lacked a concept of a holding. In the move to Alma it was decided that these article level records, monographic as they were (i.e. they are records for single articles) should sit at monographic holdings, separating them further from the serial records to which they belong. For example, an article level record for an article in the Journal of Zoology might be sitting on a different holding than the journal in which it resides. Thirdly, there is a real mixed quality of citations provided in the article records. These citations are often abbreviated, could use journal titles which are not the main entries in our catalogue and thus cannot be cross referenced easily, and may be in MARC21 fields not actually displayed in Primo VE. Finally, these records were not necessarily correctly

Bibliographic

022	— a 0952-8369
022	— a 1469-7998 (Internet)
035	— a (UK-LoNHM)11310-44nhm_inst
035	— a Catkey 11948
035	— a (Sirsi) 11310
040	— a UK-LoNHM b eng e rda
210	00 a J. Zool. Lond.
245	00 a Journal of Zoology : b Proceedings of the Zoological Society of London.
246	10 a Proceedings of the Zoological Society of London
246	13 a Journal of Zoology (London, England: 1987)
264	_1 a London : b Published for the Zoological Society of London by Academic

Holdings (1 - 5 of 5)

Sort by: **Alma Ranking** ▼

1	South Kensington / Zoology Serials Call Number SERIALS S 1A	Holdings ID 22145389340002081 ● Vol. 146 (1965) - Vol.315 (2021)
2	Tring / Tring Serials Call Number SERIALS S 130 A	Holdings ID 22162403190002081 ● Vol. 146-236 (1965-1995)
3	South Kensington / Entomology Serials Wandsworth Call Number SERIALS S 18	Holdings ID 22162149440002081 ● Vol. 146-180 (1965-1976)
4	South Kensington / Zoology Discard Call Number MAMMALS 17.1	Holdings ID 22160453480002081 ● 1 items out of 1 available
5	South Kensington / Zoology Mollusca Call Number MOLLUSCA S 8	Holdings ID 22226973490002081 ● 1 items out of 1 available

Holding

Item

	Barcode	Library	Location	Item Call Number	Call Number	Volume	Description	Year	ora ry Loc ati on	Status	Proces s type	Acce ss Num ber	Receivin g date
<input type="checkbox"/>	000399958	South Kensi...	Zoology Serials	-	SERIALS S 1A	Vol.315	Vol.315 No.4 (2021 Dec)	2...	No	Item in place	-	-	05/01/...
<input type="checkbox"/>	000337976	South Kensi...	Zoology Serials	-	SERIALS S 1A	Vol.315	Vol.315 No.3 (2021 Nov)	2...	No	Item in place	-	-	08/12/...
<input type="checkbox"/>	000353253	South Kensi...	Zoology Serials	-	SERIALS S 1A	Vol.315	Vol.315 No.2 (2021 Oct)	2...	No	Item in place	-	-	02/11/...
<input type="checkbox"/>	000354125	South Kensi...	Zoology Serials	-	SERIALS S 1A	Vol.315	Vol.315 No.1 (2021 Sep)	2...	No	Item in place	-	-	28/09/...
		South	Zoology				Vol.314			Item in			

Figure 1: Serial inventory

identified as articles by Primo VE which would confuse the user even further¹. All this has meant that, although there is a huge potential benefit to the user of having these records in the catalogue, this material is often very hard to find and obtain.

Since the move to Alma we have been addressing this problem on a record-by-record basis – utilising the MARC21 773 field to add title and, where available, barcode and bibliographic ID (MMSID) to relate the article and serial records. Alma takes the information entered in a 773 and will use this to generate holdings information on the part record. Which means that once a 773 is added with item level information it will display on the article record in Primo and hence direct the user to the serial volume they require (see [Figure 2](#)).

Details

Author: [Birkhead, Mike.](#) >

Title: Causes of mortality in the mute swan *Cygnus olor* on the River Thames
M. Birkhead.

Format: p. 15-25 : ill.

Series: [Journal of zoology 0952-8369 ; Vol. 198, 1982](#) >
[Journal of Zoology Vol. 198, 1982](#) >

Language: English

General note: 000196827

Taxonomic Name: *Cygnus olor*. Mute Swan

Locate Item

NHM staff sign-in for request options
External visitors click [here](#) for requesting items. [NHM Staff Sign in](#)

[← BACK TO LOCATIONS](#)

LOCATION ITEMS

South Kensington

→ RELATED TITLE: *Journal of Zoology: Proceedings of the Zoological Society of London.*

May be available, Zoology Serials; SERIALS S 1A

Note: Binding for Vols 211-215 incorrectly indicates Series A

Holdings: Vol. 146 (1965) - Vol.315 (2021)

Item in place (0 requests) Vol.198 (1982)

Loanable Item Barcode: 000196827

Figure 2: Example of a linked record

This link also means that if this serial is moved or another enumeration is updated, it won't affect the link, meaning a dynamic connection can be maintained. This is an elegant solution and, in this system, has discovery benefits. However, it is quite labour intensive as it requires work across item, holding and bibliographic level records. We

¹ This is down to a quirk of how Alma and PrimoVE recognise articles and a misreading of the MARC21 guidance on the a vs. b indicator in LDR/07. The table detailing how the Primo VE article resource type is mapped from MARC21 LDR/07 bibliographic level indicates that the only value that will generate a resource type of article in Primo VE is "b" – serial component part ([Ex Libris, no date b](#)). However, this refers to a serialised component such as an editorial in a newspaper rather than a component part of serial. The value of "a" – monographic component part – should also be a valid mapping to the article resource type and indeed the vast majority of articles in journals are themselves monographic, so would be properly encoded as LRD/07 = "a" ([Library of Congress, 2016](#)).

have calculated that it takes approximately 10 minutes to complete each link. As there are approximately 100,000 links to make, simply working through them manually would take over 5 years of work which is clearly not ideal!

To try and work through this issue we have established some semi-automated approaches using the functionality of Alma to delete items and holdings en masse. However, the application of linking the child record to the parent's based on the citation has always proved elusive. It was this that lead us to thinking about some more computational approaches and to contact the NHM's AI team.

Working with the AI Team

We developed a two-stage pipeline to link child records with their corresponding parent records (e.g., journals or series) using natural language processing and record linkage techniques ([de Bruin, 2023](#)). The goal was to match child items with appropriate parent candidates based on title similarity and associated metadata, including information such as year of publication, volume, and issue number.

Stage 1: Metadata Extraction and Normalisation

The initial step involved normalising and extracting structured metadata from both parent and child datasets. In the parent dataset, metadata was embedded within a free-text description field. This field was parsed using a rule-based NLP pipeline built on the spaCy library ([Honnibal, Montani, Van Landegem and Boyd, 2024](#)) specifically employing the Matcher component to identify patterns corresponding to metadata elements such as volume numbers, publication years, and issue identifiers.

In the child dataset, titles and the 773 field were similarly parsed. The series titles were processed using the same spaCy-based rule system, while additional metadata embedded in the 773 field was extracted using regular expressions. All titles were normalised (e.g., lowercasing, punctuation stripping), and any metadata extracted was retained as key-value pairs for downstream matching.

Stage 2: Candidate Matching and Best-Fit Selection

The second stage focused on identifying the most likely parent record for each child item. This was carried out using the RecordLinkage library, which provides tools for probabilistic and deterministic record linkage ([de Bruin, 2023](#)).

We applied a hybrid fuzzy matching strategy that leveraged both Jaro-Winkler similarity and partial ratio string matching ([Bachmann, 2024](#)). For each child title, candidate parent records were selected if they had either:

- The highest Jaro-Winkler similarity to the child title, or
- The highest partial ratio similarity, provided their Jaro-Winkler score was within 0.2 of the top match.

This ensured tolerance for partial or noisy matches without sacrificing precision.

Among the shortlisted candidates, a deterministic matching step compared each child record's structured metadata against that of the parent candidates. This comparison was conducted using a strict multi-field match: any candidate for which a field (e.g., volume, year) did not match the corresponding child field was excluded. If multiple candidates passed this filter, the candidate with the highest number of matching fields was selected as the best match.

Evaluation and Performance

To assess the performance of the parent-child linking system, a manually curated gold-standard dataset was assembled by the library team. This dataset comprised approximately 3,000 child records that had been manually linked to their corresponding parent journal entries. The pipeline achieved an overall accuracy of 68% when evaluated against this gold-standard dataset. Accuracy was defined as the proportion of child records for which the system's predicted parent matched the manually assigned parent. This result suggests that the title-based matching combined with deterministic metadata comparison performs reasonably well in structured contexts.

When applied to the full, real-world dataset — for which no gold-standard manual labels were available — the system yielded predictions for approximately 60% of child records. Analysis of the remaining 40% of unmatched cases revealed that the series title in the child record did not sufficiently resemble any parent title, resulting in no candidates being passed forward to the deterministic metadata comparison stage. This points to limitations in the fuzzy matching threshold and strategy, which may fail to accommodate certain types of string variation, abbreviation, or inconsistency.

Despite these limitations and given the amount of records that we need to work with, these matching percentages are excellent and we are currently working on ways to try and improve the matching with some human tweaking and intervention. Once the model has output its predicted links, we can update our records utilising a combination of MarcEdit and OpenRefine to add 773 fields to each child record, linking it to its respective parent item. We can then run jobs in Alma to remove the item and holdings records associated with these child records. What we are left with are article records linked both in the bibliographic record but also through inventory, thus allowing our users to find and request the precise volumes they need to conduct their research (see [Figure 2](#) above).

Conclusion

This paper has demonstrated a practical and effective approach to addressing a long-standing challenge in the NHM's library catalogue: linking of article-level records to their parent journal holdings. The traditional manual process of creating these vital connections was prohibitively time-consuming, posing a significant barrier to user

access and efficient resource management. By developing a two-stage pipeline leveraging natural language processing and advanced record linkage techniques, we have successfully semi-automated this complex task.

This project not only enhances the discoverability and obtainability of valuable research materials for our users but also provides a dynamic, resilient linking mechanism that adapts to changes in physical holdings. More broadly, this work exemplifies how computational methods can empower metadata professionals to tackle large-scale, repetitive tasks with limited resources, freeing up human expertise for more nuanced and strategic cataloguing efforts. We believe this methodology holds significant potential for extension to other types of complex record relationships, ultimately fostering a more accessible and user-centric library catalogue.

References

- Bachmann, M. (2024) *rapidfuzz/RapidFuzz* (v3.11.0) [Software]. Zenodo. Available at: <https://doi.org/10.5281/zenodo.14509091> [Accessed: 31 May 2025]
- Breeding, M. (2007) 'Next-Generation Library Catalogs: Chapter 1 Introduction', *Library Technology Reports*, 43(4). Available at: <https://librarytechnology.org/document/18344/> [Accessed: 31 May 2025]
- British Film Institute (no date) *Simple Search*. Available at: <https://collections-search.bfi.org.uk/web/search/simple> [Accessed: 31 May 2025]
- de Bruin, J. (2023) *Python Record Linkage Toolkit: A toolkit for record linkage and duplicate detection in Python* (v0.16) [Software]. Zenodo. Available at: <https://doi.org/10.5281/zenodo.8169000> [Accessed: 31 May 2025]
- Ex Libris (no date a) *Introduction to Alma Inventory*. Available at: [https://knowledge.exlibrisgroup.com/Alma/Product_Documentation/010Alma_Online_Help_\(English\)/040Resource_Management/050Inventory/010Introduction_to_Alma_Inventory#](https://knowledge.exlibrisgroup.com/Alma/Product_Documentation/010Alma_Online_Help_(English)/040Resource_Management/050Inventory/010Introduction_to_Alma_Inventory#) [Accessed: 31 May 2025]
- Ex Libris (no date b) Mapping to the Display, Facets, and Search Sections in the Primo VE Record. Available at: [https://knowledge.exlibrisgroup.com/Primo/Product_Documentation/020Primo_VE/Primo_VE_\(English\)/120Other_Configurations/Mapping_to_the_Display%2C_Facets%2C_and_Search_Sections_in_the_Primo_VE_Record#MARC21_and_KORMARC_Resource_Type_Mapping](https://knowledge.exlibrisgroup.com/Primo/Product_Documentation/020Primo_VE/Primo_VE_(English)/120Other_Configurations/Mapping_to_the_Display%2C_Facets%2C_and_Search_Sections_in_the_Primo_VE_Record#MARC21_and_KORMARC_Resource_Type_Mapping) [Accessed: 31 May 2025]
- Honnibal, M., Montani, I., Van Landeghem, S., and Boyd, A. (2024) *spaCy: Industrial-strength Natural Language Processing in Python* (v3.7.5) [Software]. Zenodo. Available at: <https://doi.org/10.5281/zenodo.1212303> [Accessed: 31 May 2025]
- Library of Congress (2016) *Leader (NR)*. Available at: <https://www.loc.gov/marc/bibliographic/bdleader.html> [Accessed: 31 May 2025]
- Riva, P., Le Bœuf, P., Žumer, M. (2024) *IFLA Library Reference Model: A Conceptual Model for Bibliographic Information*. International Federation of Library Associations and Institutions (IFLA). Available at: <https://repository.ifla.org/handle/20.500.14598/40.2> [Accessed: 31 May 2025]

Sonawane, C.S. (2017) 'Library Discovery System: An Integrated Approach to Resource Discovery', *Informatics Studies*, 4(3), 27-38.

Augmenting cataloguers

planning an AI Agent to generate MARC21 records

Sheldon Korpet

Library Systems and Discovery Advisor, Manchester Metropolitan University Library

Nathalie Rees

Library Systems and Discovery Librarian, Manchester Metropolitan University Library

Received: 02 June 2025 | Published: 17 June 2025

ABSTRACT

This article outlines the planning and development of an AI-powered agent designed to assist with generating records at Manchester Metropolitan University's Library and Cultural Services. The team explores the use of Large Language Models (LLMs) and Retrieval Augmented Generation (RAG) to automate and enhance cataloguing processes, particularly for unique and multilingual materials in the Special Collections Museum and Manchester Poetry Library.

The article discusses technical challenges, ethical considerations, and the importance of maintaining human oversight to ensure quality and transparency. It also details the system architecture and outlines future goals such as multilingual support and automated record enhancement. This project not only aims to improve efficiency but also empowers staff through upskilling and deeper engagement with emerging AI technologies.

KEYWORDS metadata enhancement; cataloguing automation; AI; Large Language Models; Retrieval Augmented Generation

CONTACT Sheldon Korpet ✉ S.Korpet@mmu.ac.uk 🏠 Manchester Metropolitan University
Nathalie Rees ✉ N.Rees@mmu.ac.uk 🏠 Manchester Metropolitan University

Introduction

We are part of the Digital Library Team operating within the Library and Cultural Services Department at Manchester Metropolitan University. Operating in a team of six colleagues, we use WorldCat (OCLC) as our Library Management System supplier. Our department has never had a dedicated full-time cataloguer and there is no funding for a specialist role.

We catalogue unique and self-published items for our Special Collections Museum and multilingual Manchester Poetry Library. However, cataloguing these items is a time-intensive process and there are competing demands for our expertise.

Our aim is to build an AI-powered Agent to assist in the cataloguing process. Albada (2025) defines agents as a piece of software, utilising a Large Language Model, to undertake a task without explicit instructions to initialise each individual action. The system will use an LLM to generate MARC21 format records using RDA standards and FAST subject headings. The intent being to improve cataloguing efficiency, reduce administrative load and improve item processing times.

While we acknowledge there have been significant AI advancements since 2022, cataloguers will continue to oversee the workflow. Ultimately, this AI system will suit our library staffing context, freeing staff to undertake innovative digital initiatives, such as horizon scanning, programming, and user experience research.

Past Automation Attempts

Previous in-house attempts have been made to automate the cataloguing process with various technologies (OCR, FAST API tools, document information extraction). However, an end-to-end system has not been implemented before.

Historically, AI has been rules-based but with the introduction of Large Language Models (LLMs) like transformers, it means AI can be flexible and competently dealing with natural language input (Gonuguntla, Meghana and Prajwal et al., 2024). We have experimented with commercial Generative Artificial Intelligence (GAI) tools like Copilot and DeepSeek. However, hallucinations are not uncommon, and staff found the results were mixed. Our conclusions were similar to industry findings from OCLC (Urban, 2024) as a cataloguer needs to read all the metadata generated to spot the mistakes as well as ensure it is correct and consistently conforms to our cataloguing processes.

GAI has proven itself capable of creating basic records, especially for popular titles, which can then be improved by cataloguers. However, generating DDC or subject headings is not always accurate. Creating successful prompts to produce useful records can be more time consuming than creating the record from scratch. The LLMs can bring in incorrect book information despite being given ISBNs or links to publisher pages.

Our aim is to tackle the unique, historic and self-published items we have purchased or received as donations to our Special Collections and Poetry Library. Therefore, we are looking to create a system capable of handling these items, which may not have standard layouts.

Discussion – Concerns, Barriers and Ethics

We will prioritise the accuracy and quality of records produced, mitigating the risk of introducing poor-quality records into WorldCat. However, this conflicts with internal pressure to get items out onto the shelves quickly, particularly if there is a large acquisition order or donation. There is a constant dilemma of whether it is better to create complete records in full or get a basic record on our system so items can be

discovered quickly but has to be revisited and improved later. GAI allows us to outsource speed by instantly drafting a record, which then allows staff to focus solely on accuracy, rather than searching for details and transcription.

Practical technical barriers to undertaking this activity include:

- staff skills and training
 - accessing funding and training resources,
 - having the time to learn,
 - keep up with advanced IT developments
 - building/configuring systems
- AI is a black box
 - transparency needs to be built into the systems
 - All models have some bias “baked in” due to how they are created ([Resnik, 2024](#))
- making machine-readable documents
 - transcribing cataloguing processes and procedures in detail is time consuming
 - handwritten notes and non-OCR content is a challenge for machines to ingest
- I.T. infrastructure
 - having a laptop with enough computing power or access to a virtual machine
 - getting software purchase approval for pay-for-use tools and online platforms
 - installing software on institutionally managed laptops

Despite the technical barriers, by building a system ourselves we have control over the decision making and execution process. A criticism of the recent AI tools from publishers includes that the model displays overly broad sensitivity ([Tay, 2025](#)). By building a custom system, we can control the configuration to meet our requirements. Should we find the model is too focused or too broad, we can easily alter the parameters of the model.

Additionally, there is a chance OCLC may release a similar Agent to compete with the Alma offering that other institutions are currently testing. This would duplicate our efforts, but it's something that we'd like to conceive rather than wait for the market to catch up. The main reason is because we can build our values into the Agent. We would prioritise the system's transparency and observability, allowing our cataloguers to understand how it is operating internally, thus improving usability and enabling user insights. This is something suppliers may not have time or inclination to explain.

We feel quite strongly about the ‘human in the loop’ approach. Cataloguing is not a science – anomalies arise, and decisions must be made to standardise record creation across staff and time. By having a human acting as a quality assessor who feeds back into the system, we aim to achieve two things:

- Consistency of output - catch errors at the source and ensure a level of standardisation across our catalogue.
- Training the agent – by having a human review and correcting AI generated records we are giving the system a high-quality dataset, allowing it a chance to see high quality records and ‘self-learn’ by embedding these into the vector database with domain knowledge

There are fears that AI could take work from humans. While we don’t have full-time cataloguers, we believe the role of cataloguer won’t disappear but its nature will change. We will still need to retain staff to oversee the operation of the system. Records need to be checked, new standards implemented, books physically labelled, new training documents written and processed into the system’s knowledge base. We should also monitor the functioning of the system. If anything, there is a need for change management and upskilling for cataloguers, rather than a risk of deprofessionalisation.

Even if we don’t take this system into our cataloguing workflow, it may be decided we should evaluate or purchase a similar product when it comes to market. The team will already possess the knowledge to critically evaluate technologies underpinned by LLMs, having had to think through and attempt the design process. Additionally, we also find projects like this often inspire other projects which would have otherwise never been discussed.

Technical Systems

Here we will describe two AI Agents. The first was a prototype and the second, design plans for a more advanced system with domain information but having access to more tools and memory.

Prototype build

We used RAG to give the LLM access to a domain-specific knowledge base. By allowing the LLM to access more information, it can aid it to generate an accurate and specific answer outside the areas increased in the model training data set ([Mendelevitch and Bao, 2025](#)). We used RAG to allow the agent to access specific information about how we catalogue for our library. LangChain was used as the agent framework, Chroma was used as the Vector database and Streamlit as the user interface. The system ran in the user's browser.

LangChain gave the document chatbot context, or a role to play. It instructs it to be a cataloguer, creating MARC21 format records, using RDA and FAST subject headings.

The documents in the vector database detail all our institutional guidance, as well as details about MARC and RDA which human cataloguers use (see 'Information sources').

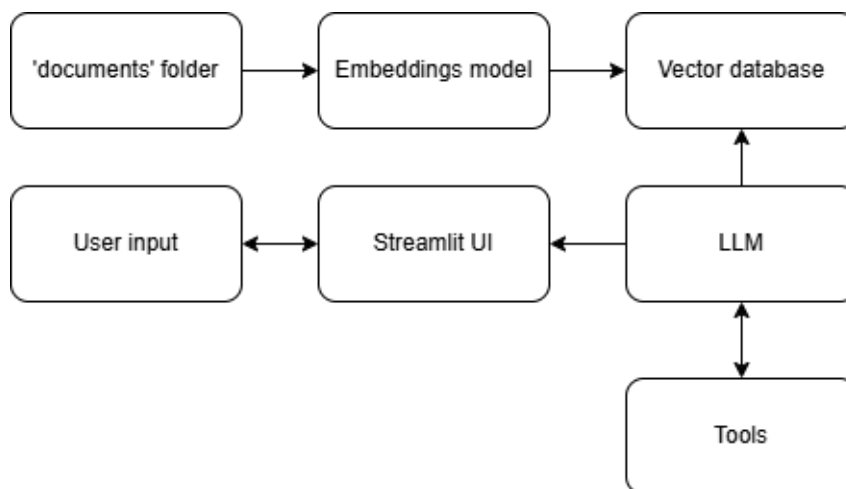


Figure 1: Simple agent system diagram

While some of the record generated was incorrect (subject headings were formatted using LCSH not FAST), the record generated showed enough improvement compared to commercial LLMs to decide to explore this technology further. We were impressed with the use of fields and noticed that the record was generating fields we often didn't have time to fill (505 for chapter titles and 520 for summary descriptions, if the information was provided).

Information sources

When deciding what to include within the RAG document store for the agent, we considered what a new member of the cataloguing team would need to know. Not just how to catalogue but also what our organisational practices and policies were. We provided the embeddings model with our policy and guidance documents, information about the DDC system, MARC21 fields, RDA documents, FAST Subject Headings and a few example records for specific collections.

Planning the new AI Agent

This is often done in the form of a 'Design Document' as best practice (Ubl, 2020). This helps share the idea, make transparent the system design plan, capture institutional knowledge and think through issues.

Design philosophy

There are five levels of complexity in Agents (Huang, 2024). Our first attempt was Level 3 – with the RAG system LLM having access to domain knowledge through the knowledgebase and basic memory and a simple web search tool. This time we will aim for Level 4 – the agent will have a cataloguing knowledge base still, but we will also build short term memory for individual conversations and a long-term memory

database. The system will be able to 'learn' in a simplistic way by saving any records generated as examples for the future.

Goals and non-goals

These aims help define the scope of the system and outline clearly how the system functions.

The system will not:

- be able to transcribe information with 100% accuracy
- be able to guarantee an accurate, complete record is generated
- be able to operate without human oversight

The system will:

- draft records almost instantly, ready for human review and correction
- have access to expert-identified and curated domain knowledge
- use tools to suggest a shelf mark
- use tools to suggest FAST subject headings
- have observability built in so threads and decision making can be reviewed by a human
- use a database for conversational persistency (long term memory)
- learn from record generation process by integrating records into its knowledgebase

System description

This agent will be given a Context on start-up (this defines its purpose):

"You are a cataloguer for an academic library, special collections museum and poetry library. You create MARC21 format records using RDA and use FAST subject headings."

It will also have access to tools and built-in memory. LangGraph works by connecting nodes; information flows through these, allowing the process to iterate somewhat. Each agent can follow the conversation as a thread, and each thread is observable using LangSmith (another Python Library). That allows the logic process to become somewhat transparent to the user if further investigation is required.

The system will use a Streamlit user interface (UI) where cataloguers can input natural language queries alongside scanned images of book materials (covers, bibliographic information, author pages etc). The images submitted to the UI will be processed via a Vision Language Model (VLM). This will extract text from the images and show the user the transcription before records are generated. Cataloguers will be able to review the transcription of this material produced by a VLM, view and re-submit

records with errors. The system will be designed to work locally on our machines and we'll utilise GitHub to hold our code repository.

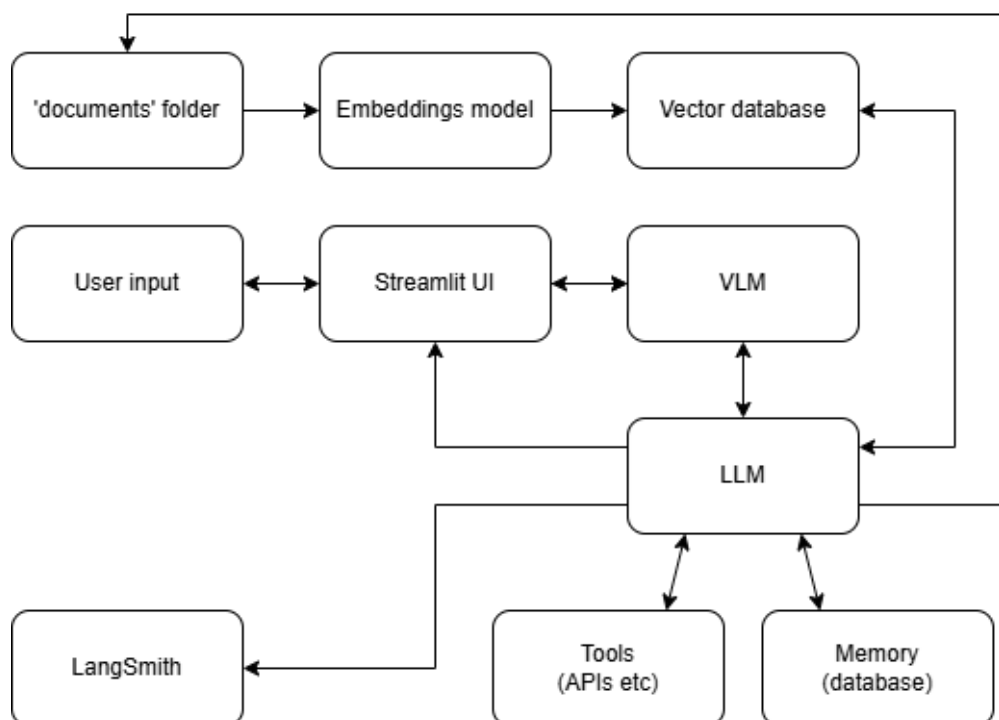


Figure 2: Advanced agent system diagram

The knowledge base will be created by a preprocessing script. An embeddings model will be utilised to make a Chroma database. This is a vector database with numeric representation of the domain documents provided to it. This allows the LLM to access the knowledge when generating a response to the cataloguer's inputs.

We plan to make use of OCLC APIs – Fast Search, Metadata. These will allow the agent to access data outside its models to choose subject heading information as well as check the system for existing records. Should the API connection fail, the system will alert the cataloguer and proceed to generate the record. The system will have clear and helpful error messages so the (likely non-programming) cataloguer will be able to understand what failed and feedback to the developer.

After the record has been generated, a copy will be saved locally as a .txt file. It will then be possible to run the preprocessing script before the next cataloguing session and give the agent an expanding bank of human quality assured cataloguing examples. As we are not finetuning a model, we can swap these records at will – should we want to create only a certain type of record e.g. we're working on a zine project, it should be possible to focus the agent on a specific subject to catalogue.

Framework

LangGraph is a Python framework that expands possibilities for more complex systems. It works by connecting nodes. Information flows through these, allowing the process to iterate somewhat. Functions include tool calling (allowing the LLM to grab data from APIs), built-in memory (both within the current conversation or thread and between sessions), and the use of an interrupt function (which stops the chatbot to allow the human to intervene and increase output quality).

Each agent is able to follow the conversation as a thread, and each thread is observable using LangSmith (another Python Library). That allows the logic process to become somewhat transparent to the user if further investigation is required. This is necessary as the developer gives the LLM access to tools – it doesn't specify when and which tools it should call or use.

The reasons for choosing LangGraph include:

- LangChain, the technology we previously used for RAG, is depreciating some of its functions to LangGraph.
- LangGraph is also popular, and this means there are plenty of online courses and community support.
- The documentation is comprehensive and well-defined.
- Additionally, it's built to utilise LangSmith, which is a platform which allows us to observe the application and evaluate performance.
- These are industry standard tools, so while unnecessary for a small project hosted remotely from our local machines, it's good practice and will provide a project to enable skills acquisition in the team.
- The framework is scalable, and it will give us options in the future should we wish to expand into multi-system agents (Level 5 - Agentic AI – see [Huang, 2024](#)).

Model choices

This is a difficult task as there are many model options and many evaluation tests. The key thing is to choose a model trained and capable of doing the task you are asking it to do ([Kirmer, 2024](#)).

It is possible to use commercial models e.g. Anthropic's Claude Sonnet model or ChatGPT 4.x. However, we would prefer to choose an open-source model as it's free. One tool which can be useful for identifying high performing LLM models are Hugging Face Leaderboards ([Hugging Face, 2025](#)).

We will likely need three models for the system:

- a Vision Language Model (VLM) - to transcribe text from item scans
- an LLM to create the record and “converse” with the user in natural language
- an Embeddings model to be used with the document store and create the vector database (domain specific knowledge which will be accessible to the LLM)

Testing Record Generation

Assessing the effectiveness of the system can be done by providing a human with the item to be catalogued and comparing the record quality (presence of fields, relevant subject headings etc) of both. The system will be able to produce a record faster every time, so what we will be directly comparing/measuring is how accurate the transcription is, how well the record represents the item, and how much detail the records have.

The system can also be compared directly to other commercial LLM services without domain knowledge via RAG. This should demonstrate how having access to process documents and specialist knowledge changes generative record production.

Future Development

We see potential for development in four areas:

- Collaboration - we’re discussing with other libraries whether projects like this can be made into collaborative, open-source efforts.
- Technology - We would like to explore using this technology to automatically enhance current records, such as adding more subject headings to aid discovery as well as evaluating its ability to catalogue in other languages for our Poetry Library.
- Skills and Training – projects like this allow a framework for staff to develop, upskill, engage in interesting and challenging work.
- Efficiency – reducing administrative tasks, like typing and transcribing, and creating time efficiencies.

Glossary

Agent - a piece of software, utilising a Large Language Model, to undertake a task, make decisions without humans. More advanced ones have memory and call tools.

Hallucination - mistakes the Generative AI model makes when creating an answer. Often a result of poor training data, incorrect assumptions, or biases ([Google, 2024](#))

Retrieval Augmented Generation (RAG) - the improvement of LLM outputs by allowing them to access domain specific information.

LangGraph - a Python library. It is prewritten code which can create the framework for generative AI agents without coding everything from scratch (Taulli, 2024).

Large Language Models (LLM) - Machine Learning models, often with billions of parameters. They are trained on massive datasets and often generate language or perform natural language processing.

Machine Learning – computers learn and identify complex patterns from data without being explicitly programmed by humans.

Python – a versatile programming language used mostly for AI, Machine Learning, automation.

Vision Language Models (VLM) - Large Language Models which can process images, extract key data like colours and shapes, then make observations available to other models via mathematical representations.

References

- Albada, Michael (2025) *Building Applications with Agents*. O'Reilly. Available at: <https://learning.oreilly.com/library/view/building-applications-with/9781098176495/> [Accessed: 20 May 2025]
- Gonuguntla, H., Meghana, K., Prajwal, T.S., Koundinya, N.V.S.S., Sai, K.N. and Jain, C. (2024) 'Evaluating Modern Information Extraction Techniques for Complex Document Structures', *International Conference on Electronic, Computer and Energy Technologies (ICECET 2024)*, Sydney, 25-27 July 2024. Available at: <https://doi.org/10.1109/ICECET61485.2024.10698618> [Accessed: 30 April 2025]
- Google (2024) *What are AI hallucinations?* Available at: <https://cloud.google.com/discover/what-are-ai-hallucinations> [Accessed: 28 May 2025]
- Huang, Yu (2024) *Levels of AI Agents: from Rules to Large Language Models*. Available at: <https://doi.org/10.48550/arXiv.2405.06643> [Accessed: 28 May 2025]
- Hugging Face (2025) *Leaderboards and Evaluations*. Available at: <https://huggingface.co/docs/leaderboards/en/index> [Accessed: 28 May 2025]
- Kirmer, Stephanie (2024) *Choosing and Implementing Hugging Face Models*. Available at: <https://www.stephaniekirmer.com/writing/choosingandimplementinghuggingfacemodels> [Accessed: 28 May 2025]
- Langchain (2025) *Langgraph*. Available at: <https://langchain-ai.github.io/langgraph/> [Accessed: 19 May 2025]
- Mendelevitch, Ofer and Bao, Forrest (2025) 'Rag vs. fine-tuning' in O. Mendelevitch and F. Bao *Hands-On RAG for Production*. Available at: <https://learning.oreilly.com/library/view/hands-on-rag-for/9798341621701> [Accessed: 20 May 2025]

- Resnik, Philip (2024) *Large Language Models are Biased Because They Are Large Language Models*. Available at: <https://doi.org/10.48550/arXiv.2406.13138> [Accessed: 23 May 2025]
- Taulli, Tom (2024) *What is LangGraph?* O'Reilly. Available at: <https://learning.oreilly.com/videos/what-is-langgraph/0642572071776/> [Accessed: 20 May 2025]
- Tay, Aaron (2025) 'Deep Dive into Three AI Academic Search Tools', *Katina*, (20 May). Available at: <https://doi.org/10.1146/katina-052025-2> [Accessed: 20 May 2025]
- Ubl, Malte (2020) 'Design Docs at Google', *Industrial Empathy*, 6 July. Available at: <https://www.industrialempathy.com/posts/design-docs-at-google/> [Accessed: 28 May 2025]
- Urban, Richard (2024) 'Keeping up with next-generation metadata in archives and special collections', *Hanging Together*, 17 December. Available at: <https://hangingtogether.org/keeping-up-with-next-generation-metadata-in-archives-and-special-collections/> [Accessed: 8 May 2025]

Smart enough to mislead

the functional shortcomings and ethical dilemmas of generative AI use in metadata work

Fran Frenzel  0009-0009-3890-5335

Metadata Analyst, London School of Economics and Political Science

Received: 10 June 2025 | Published: 17 June 2025

ABSTRACT

This article critically examines the applicability of generative AI in library metadata creation and cataloguing, arguing that despite growing interest and experimentation, such technologies remain fundamentally unsuited for this domain. Drawing on recent literature, surveys, and institutional case studies, the author demonstrates that generative AI tools consistently produce metadata outputs that are unreliable, inconsistent, and ethically problematic. While machine learning offers potential in specific, supervised metadata functions, generative AI's reliance on probabilistic outputs, lack of transparency, and tendency to hallucinate undermine the accuracy and reliability essential to cataloguing. The article also explores the broader ethical implications of AI adoption in libraries, including issues of bias, environmental impact, copyright concerns, and labour exploitation. The author argues that fully automated metadata creation using generative AI is neither technically viable nor ethically responsible and instead advocates for cautious, critically informed AI integration, emphasising the continued necessity of human oversight and ethical scrutiny in metadata work.

KEYWORDS generative AI; metadata creation; metadata enhancement; AI ethics

CONTACT Fran Frenzel  f.frenzel@lse.ac.uk  London School of Economics and Political Science

In early autumn 2023, *Information Technology and Libraries* published an article proclaiming that ChatGPT can “generate accurate MARC records using RDA and other standards such as the Dublin Core Metadata Element Set” (Brzustowicz, 2023, p. 2). I was intrigued and began reading, quickly followed by bewilderment and dismay about what was labelled “an accurate and effective record” (Brzustowicz, 2023, p. 2) and “comparable to the professional catalogers’ work” (Brzustowicz, 2023, p. 3), but was riddled with errors and hallucinations¹. All the article demonstrated to me was ChatGPT’s utter inadequacy for cataloguing, and I was horrified by the potential consequences of assertions, such as the ones quoted above, being taken at face value. On the positive side though, it got me interested in AI in relation to library metadata work. So, nearly 2 years on, are we any closer to AI cataloguers?

¹ The relevant mailing lists were not amused either and some responses were published in the next issue of *Information Technology and Libraries* that outline the problems very well (DeZelar-Tiedman, 2023; Amram, Malamud & Hollingsworth (2023) and Floyd (2023)).

Yes and no. In terms of AI technologies being used for cataloguing and record enhancement, yes, there are many possibilities to assist cataloguers or streamline processes. But in terms of generative AI being the AI technology to do the job, very much not. Why? Well, generative AI is singularly concerned with creating the most likely response to a prompt; whether this most likely response is factually correct or not is of no importance to it. This is in direct opposition to what one is doing and what is important when cataloguing: truthfully representing the resource, not representing it with the statistically most likely information. That is not to say that there aren't parts of records or tasks in metadata creation where generative AI's statistical approach might be useful, such as summary creation, for example.

Throughout this article I make a distinction between

- **AI**, meaning AI overall,
- **machine learning**, meaning the field of study with AI research, and
- **generative AI**, meaning a specific subtype of AI utilising generative machine learning models, deep learning and neural networks.

To understand how generative AI functions, it is useful to get a grounding on how large language models (LLMs) work. The Financial Times has published a very accessible primer on how LLMs function and the enormous difference the development of the transformer deep learning architecture has made, leading to generative pre-trained transformers (GPT, a type of large language model) which underpin generative AI tools such as OpenAI's ChatGPT, Microsoft's Copilot, Google's Gemini, Anthropic's Claude or Meta's Llama ([Murgia et al., 2023](#)). Hicks, Humphries and Slater ([2024](#)) also give a good and accessible explanation of LLMs' functionality. With this information in mind, it becomes easy to understand why LLMs produce superficially convincing-looking records that, upon closer inspection, reveal a multitude of problems.

Opportunities

Chen and Li ([2024](#)) published results of a survey conducted in early 2024 to gauge cataloguing and metadata professionals' perception of AI in relation to their roles. They found that AI does not yet play a significant role in respondents' jobs ([Chen and Li, 2024](#), p. 321), but that there is "a growing belief [sic] that AI would play an increasingly important role in their work" ([Chen and Li, 2024](#), p. 322). Among the questions asked were in which areas of respondents' work AI is currently used, and in which areas of metadata creation they think AI would be most beneficial. Translation and summary creation were mentioned most in relation to current use and also ranked highest in areas most benefiting, followed by subject headings and classmarks. I am unsure what was meant by "physical description" and "creators/contributors" in the survey, but as data transcription was not mentioned yet, my assumption is that it is covered by these two categories and/or "other" ([Chen and Li, 2024](#), pp. 322-323).

These survey results tally with my own experience in AI use for cataloguing, and where I believe AI can be of most benefit to cataloguers.

While born digital resources and digitised resources might spring to mind as the most likely candidates for AI-assisted metadata work, Lowagie (2024) has shown that there are also possibilities for physical resources.

Beyond individual record creation, machine learning, I think, has huge potential to improve metadata management tasks, though I cannot see generative AI to be useful in this respect. Data management tasks need to be transparent and produce repeatable and consistent results, all of which generative AI outputs certainly are not.

OCLC, for example, are using machine learning (not generative AI) to deduplicate WorldCat records and have removed millions of duplicate records from WorldCat with this approach (OCLC, 2025). Another great example of machine learning assisting in bulk tasks is Cornish and Scott's article in this number of Catalogue and Index (Cornish and Scott, 2025).

Interestingly, the Chen and Li's survey also found that "most respondents didn't find AI tools had significant help of [sic] either quality or efficiency of their cataloging work" (Chen and Li, 2024, p. 324). I believe this is grounded in the problems with accuracy and reliability of AI outputs. If one needs to double-check everything, there is no significant time saving or efficiency in the use. In my experience, thoroughly checking and, if needed, amending a record takes the same or even more time than creating it in the first place. At LSE we investigated ExLibris' generative AI enriched Community Zone records (ExLibris, no date a; ExLibris, no date b; York and Hanegbi, 2024) some months ago and found the assigned subject headings often on the too broad side, while the quality of summaries varied significantly depending on the type of publication. We also trialled using generative AI for JEL code² assignment. In addition to not seeing a time saving due to the need to check the AI output, colleagues reported that, while they found it to be a useful assistive tool, they also felt their ability to assign codes without the use of AI and familiarity with the vocabulary as a whole declined. A sentiment that is also echoed in Chen and Li's survey results as "worry about over-reliance" (Chen and Li, 2024, p. 324)

What can machine learning and generative AI do?

Subject headings and classification are probably the areas that have been looked at the most for automation so far.

The National Library of Finland's well-known Annif³ tool (Suominen et al., 2023) has been around since 2019 (National Library of Finland, 2025) and is in use in various libraries in either fully automated workflows or human-supervised ones (Inkinen,

² The JEL (*Journal of Economic Literature*) classification is a commonly used classification scheme for scholarly works in economics. <https://www.aeaweb.org/econlit/jelCodes.php?view=jel>

³ Annif is an example of the application of machine learning and thus AI, but it does not, as far as I am aware and understand, make use of generative AI.

[Lehtinen and Suominen, 2025](#), p. 3). An Annif user survey, however, also reveals potential problems with implementing such a tool:

- The technical expertise needed. The survey shows that the most encountered problems with the tool are of a highly technical nature ([Inkinen, Lehtinen and Suominen, 2025](#), pp. 3-4), which suggests levels of technical expertise are needed for an implementation that most libraries will be unable to shoulder by themselves or at all.
- The resources needed. The implementing institutions are big ones (national libraries, university libraries) ([Inkinen, Lehtinen and Suominen, 2025](#), p. 1), which suggests resourcing that will be out of reach for most others.
- The tool not delivering the expected results/time savings ([Inkinen, Lehtinen and Suominen, 2025](#), p. 4).

The German National Library (DNB) is using Annif for its “Erschließungsmaschine” (EMa, “subject cataloguing machine”). Results published in 2021 regarding the assignment of subject headings (GND⁴ descriptors) showed a rather worrying 10% of assignments having been assessed by subject experts as “wrong” and 22% as “less useful” ([Uhlmann and Grote, 2021](#)). The DNB has since also started assigning DDC short numbers (a simplified classification system that the DNB developed ([Deutsche Nationalbibliothek, 2023](#))) automatically. The performance metrics for this indicate a very mixed picture as well, with some categories scoring pretty well, but others rather badly ([Poley et al., 2025](#), pp. 12-13). Poley et al. ([2025](#)) state that the volume of available training data is an important factor in the model’s performance, but also that it does not seem to be the only criterion. As other criteria are not mentioned, I assume the authors do not know either (black box).

Golub et al. ([2024](#)) have also used Annif to conduct research on automated Dewey Decimal Classification numbers, working with data in the Swedish union catalogue. They achieved a 66.82% accuracy rate on assigning three-digit DDC numbers by combining the results of four classification algorithms. During the research, they discovered that classifying fiction posed a problem and identifiable records for fiction were excluded from the set of records to be classified, which improved the accuracy rate to the figure mentioned above. However, fiction records remained in the training dataset and thus fiction headings are assigned. The accuracy rate might be improved by excluding fiction from the training dataset.

Chow, Kao and Li ([2024](#)) experimented with assigning Library of Congress Subject Headings using generative AI and found that generative AI only produced usable outputs in about half of their samples. They thus conclude that “while ChatGPT can access an internalized corpus of the LCSH and MARC 21 [sic] bibliographic records, the model struggles with validity, specificity, and exhaustivity in its generated subject headings” ([Chow, Kao and Li, 2024](#), p. 585) and that “in order to ensure accuracy and

⁴ The GND (Gemeinsame Normdatei, “integrated authority file”) is the standard German-language authority file and contains personal and corporate names as well as subject headings.

reliability of the cataloging process, the involvement of human catalogers remains an essential prerequisite" ([Chow, Kao and Li, 2024](#), p.586).

The Exploring Computational Description experiments by the Library of Congress (LoC) in cataloguing eBooks via generative AI also showed low quality scores getting nowhere near the goal set, but it showed some promise in extracting author names, titles and identifiers ([Weinryb-Grohsgal, Potter and Saccucci, 2024](#); [Saccucci and Potter, 2024 b](#), p. 6; [Library of Congress, no date](#)). As with the DNB activities, the LoC experiment also highlights the importance of the "quality and robustness of the training data" ([Library of Congress, no date](#)) for the success of AI-generated records. The overall conclusions of the LoC experiment, at the time of writing, remain:

- No current generative AI tool returns good enough results to run automatic cataloguing.
- "Human-in-the-loop" workflows are possible though, and should be explored.

The third phase of the experiment started in August 2024 ([Weinryb-Grohsgal, Potter and Saccucci, 2024](#)), but no results have been published yet.

The technical hurdle for the implementations and experiments above is rather high, but less tech-intensive solutions are possible as well, as Lowagie showed with the KBR's approach ([Lowagie, 2024](#)). He implemented an AI-driven solution to bottlenecks in metadata creation and retro-cataloguing backlogs by employing Power Apps to extract information from photographs of title pages, thus enabling cataloguers to concentrate on ensuring correct information rather than data entry. Lowagie also introduced creating custom application profiles with Power Apps and using them to validate records in the catalogue.

Another low technical hurdle experiment was undertaken by Taniguchi ([2024](#)), who used the illustrative sources of information in *Maxwell's Handbook for RDA* to generate records using ChatGPT. The conclusion here is also that the generative AI produced records with significant errors and "struggled with complex bibliographic patterns and nuanced cataloging rules", but could conceivably be used as an assistive tool for human cataloguers ([Taniguchi, 2024](#), p. 544).

All these examples show that AI technologies can be leveraged to assist in cataloguing and metadata maintenance, but, apart from Taniguchi ([2024](#)), they are also all very far from the "prompt in chat to ingestible record" scenario that started this piece. I fully agree with Moulaison-Sandy and Coble ([2024](#)) in their assessment that "the perception that [AI] is able to solve specialized problems in cataloging easily, with the click of a button, if only the right prompt is created, is problematic to perpetuate" but also that "now is the time to look to the future and to be creative, but with a sense of the full understanding of the limitations and affordances." ([Moulaison-Sandy and Coble, 2024](#), p. 382)

This sentiment is further echoed in a survey the Program for Cooperative Cataloging (PCC) ran in March 2024 to gauge current AI activities and what their impact on cataloguing and metadata work is. The first theme emerging from it is:

“The need to clearly communicate to library administrators and the broader cataloguing community that AI is not an easy fix or money saver. AI and ML [machine learning] technologies take time and careful consideration in order to be implemented effectively and must be done in concert with cataloging and metadata experts.” ([Program for Cooperative Cataloging, 2024](#), p. 2)

The newest generation of generative AI models are no longer pure LLMs but Large Multimodal Models (LMMs) able to handle not just text in- and output, but other media such as images, audio and video as well ([Wu et al., 2023](#)). Large Reasoning Models (LRMs) are developed with better reasoning and fact-checking abilities to improve performance ([Mollick, 2025 a](#)). On the other hand, Apple just released a paper that reckons this is all just an “illusion of thinking” and “that frontier [large reasoning models] face a complete accuracy collapse beyond certain complexities” ([Shojaee et al., 2025](#), p.1). Furthermore, there is evidence that newer models hallucinate more, and the developers and researchers do not understand why ([OpenAI, 2025](#), p. 4; [Chowdhury et al., 2025](#); [Metz and Weise, 2025](#)).

What could be promising though, is a combination of generative AI, computer vision tools, good old database queries (not everything needs to be generated new, sometimes just finding a good, existing record and verifying it is all that’s needed) and incorporation of local documentation. The use of the latter two, I recently learned, actually has a name and framework: Retrieval Augmented Generation (RAG). RAG retrieves information from external sources (e.g. a knowledge base, database, etc.) and uses this to augment the LLM response. By doing this, the LLM can return more accurate and contextually relevant responses ([Google, no date](#)).

So, yes, there are opportunities to employ AI technologies (and generative AI can be part of this) to create or enhance metadata, but can it be done with the needed reliability and accuracy, or can entire records be created? Absolutely not – at least not for the time being.

With generative AI bullshitting⁵ and hallucinations seemingly not going anywhere, through them being an inherent part of it⁶, I cannot see generative AI-driven solutions for metadata creation and enhancement being able to operate without oversight; we need the human in the loop to check outputs.

Time and emerging technologies may well change this view.

⁵ Strictly in the Frankfurtian sense outlined by [Hicks, Humphries and Slater, 2024](#).

⁶ “Despite our best efforts, they will always hallucinate” Amr Awadallah, formerly of Google and now CEO of a startup building AI tools for businesses, told the New York Times ([Metz and Weise, 2025](#)).

Should we just because we can?

Now that we have covered the technical possibilities, let's have a look at the ethical side of things. Berkowitz (2025) argues that libraries tend to choose quick adoption of emergent technologies (AI use in this case) "for the sake of being perceived as cutting-edge early adopters" over "a deeply methodical intent", i.e. they tend to opt for FOMO over slow-mo (Berkowitz, 2025, p. 52). He calls for libraries to champion AI ethics and concentrate on ethical scrutiny and developing policies and ethical frameworks for AI use rather than quick adoption.

AI bias

Bias can enter generative AI outputs in various ways:

- It can be present in the training data, for example, because the data is not representative, omits or obscures information. There can also be problems with inconsistent training data labelling. However, even data that is otherwise sound will still reflect structural and historical biases.
- Secondly, the training and inference algorithms may display bias or amplify biases in the training data.
- Further biases can be introduced through the evaluation of a model and the used benchmark dataset(s).
- Finally, models may be used in scenarios they were not intended for and thus produce biased, harmful outputs (Gallegos et al., 2024, p. 1107).

For a much deeper dive into AI biases, their evaluation, and techniques for bias mitigation, I refer you to Resnik (2025) and Gallegos et al. (2024).

What are the consequences of an agent with harmful biases creating metadata? Well, metadata that furthers and perpetuates those, of course. A lot of work has been done in libraries to overcome harmful language in records as well as to bring materials by and about marginalised groups out of their space of marginalisation and othering. By using generative AI to help us in creating and managing metadata, are we negating at least some of this work? As Corrado (2021, p. 402) asks: "how will [AI] satisfy the ethical concerns related to representation and identity in metadata?" As humans we can of course use our judgement and awareness of our own biases to ensure we do not perpetuate inequities or "hide" content behind overly broad or othering headings or by omission. I agree with Corrado that "it is yet to be seen how artificial intelligence will deal with this fluid space. Unless librarians and other advocates push for this, the answer may very well be that it won't." (Corrado, 2021, pp. 402-403)

Given the importance of representation in training data the German National Library (Poley et al., 2025) and Library of Congress (Library of Congress, no date) have found in their respective subject indexing and metadata creation experiments as well as the struggles with specificity Chow, Kao and Li (2024, p. 585) found, I have doubts that AI-

generated subject headings can adequately represent material on subjects underrepresented in the training data or on emerging subjects. Poley et al. in fact recognise that this is the case, stating that “subject areas where automatic subject cataloguing does not work, or does not work well, must first be intellectually indexed in order to generate training data to improve possible machine models.” ([Poley et al., 2025](#), p. 25)

AI’s black box nature further obscures things, making it very difficult indeed to both identify and address biases in outputs. How can one intervene in a “thinking” process whose workings are not understood even by the people developing them ([OpenAI, 2025](#); [Metz and Weise, 2025](#))?

Copyright

Copyright legislation is woefully behind AI development, and many questions are unanswered regarding copyright ownership of AI-generated content as well as copyright infringements in training AI models. Most generative AI training data is scraped from publicly available internet resources, but it also includes material from platforms that contain copyrighted material (e.g. the recent outcry over the Library Genesis dataset). Alex Reisner’s article on the subject is a scary read indeed ([Reisner, 2025](#)). He details that both Meta and OpenAI argue that their use of copyrighted material for training without a license falls under “fair use”. I believe the courts have yet to make a judgment on this, but the US Copyright Office certainly begs to differ ([United States Copyright Office, 2025](#); [Constantino, 2025](#)).

Ball ([2025](#)) adds that “To add insult to injury, academic publishers are now beginning to license access to their content to AI companies, some without providing academics the opportunity to opt out. This forces complicity on academics, turning their intellectual contributions into commodities for AI profit without their consent and with no remuneration for them or their institutions.” See also Battersby ([2024](#)) and Eaton ([2024](#)) on this subject.

If libraries engage in AI use, I think they should think hard about these issues and whether it is ethical to condone such practices by using tools built on them.

Environment

“There is still much we don’t know about the environmental impact of AI but some of the data we do have is concerning”. This is the statement of the Chief Digital Officer of the United Nations Environment Programme in a news piece summarising report findings ([UN Environment Programme, 2024](#)). Hugging Face, which provides a platform for sharing machine learning models and datasets, also acknowledges that “the nature and extent of AI’s effects are under-documented, ranging from its

embodied and enabled emissions to rebound effects due to its increased usage” ([Luccioni, Trevelin and Mitchell, 2024](#))⁷.

Generative AI companies are not forthcoming with accurate and complete data on the environmental footprint of their products, and available figures rely on lab-based research, such as that carried out by Luccioni and Strubell⁸ as well as “limited company reports; and data released by local governments” ([Crawford, 2024](#)). Annual reports of big tech companies show that they are not meeting their sustainability targets ([Barker, 2025](#)).

The negative environmental impact of AI splits between its electricity and water consumption, resources needed to manufacture equipment and the regurgitation of it as electronic waste ([UN Environment Programme, 2024](#)). It includes not just the training of models, but also their usage.

The training and operation of AI demands vast quantities of computational power, and the data centres housing the servers that do the work need electricity and water for cooling. Loads of figures are floating about on the energy consumption of generative AI data centres:

- Mollick ([2025 b](#)) states that there are now AI models in use that consume as much computing power for training as it takes to “[run] a modern smartphone for 634,000 years or the Apollo Guidance Computer that took humans to the moon for 79 trillion years”.⁹
- Zewe ([2025](#)) states that the energy demand of data centres in North America is estimated to have increased “from 2,688 megawatts at the end of 2022 to 5,341 megawatts at the end of 2023”, an increase “partly driven by the demands of generative AI”. He also mentions that the global data centre electricity consumption reached 460 terawatts in 2022, which makes it the “11th largest electricity consumer in the world, between the nations of Saudi Arabia (371 terawatts) and France (463 terawatts)” and that it is expected to be closer to 1,050 terawatts by 2026.
- The UN Environment Programme ([2024](#)) gives the example of Ireland, which hosts many data centres, stating that the International Energy Agency estimates that “the rise of AI could see data centres account for nearly 35 per cent of the country’s energy use by 2026”.
- Taking the research by Strubell, Ganesh and McCallum ([2020](#)) as the basis, Luccioni, Trevelin and Mitchell ([2024](#)) state that “training [an LLM with] 213 million parameters was responsible for ... [the] equivalent to the lifetime emissions of five cars, including fuel”. For comparison, OpenAI’s most recent

⁷ Luccioni is Hugging Face’s Climate Lead, Mitchell its Chief Ethics Scientist and Trevelin its Legal Counsel.

⁸ The reference is to [Strubell, Ganesh and McCallum, 2020](#), [Luccioni, Viguier and Ligozat, 2023](#) and [Luccioni, Jernite and Strubell, 2024](#).

⁹ In the correct unit of measurement that’s 10^{26} FLOPS (Floating point operations per second).

model, GPT-4o, allegedly¹⁰ uses 200 billion parameters ([Ben Abacha et al., 2025](#)). Zewe (2025) states that the electricity needed for training a model like GPT-3 was estimated to consume the equivalent of 120 average U.S. homes yearly energy consumption.

Training a model is just one part of it though: energy is consumed every time the model is used, and, with the rapid development of new models, training needs to be repeated for them frequently ([Zewe, 2025](#)).

In terms of model usage, Luccioni, Jernite and Strubell (2024) found clear differences between modalities, with image-based tasks and generation of new content using the most energy. A report prepared by Goldman Sachs found “that a ChatGPT search consumes around 6x-10x the power as a traditional Google search” ([Goldman Sachs, 2024](#), p. 13). This report also shows quite frightening predictions for power use by AI. As Luccioni, Trevelin and Mitchell state, “the growing energy demand for AI is significantly outpacing the increase in renewable energies – entailing substantial new [greenhouse gas] emissions and squeezing an already tight renewable energy market.” ([Luccioni, Trevelin and Mitchell, 2024](#))

Next on the list is water consumption: water is needed to cool the servers in the data centres, and it needs to be cooled so it can absorb the heat from the machines. Additionally, neither salt nor grey water can be used for this process as this damages the cooling systems. Approximately 30-40% of the electricity consumed by data centres is used for water cooling ([Luccioni, Trevelin and Mitchell, 2024](#)). The amount of water needed depends largely on the size of the data centre. The biggest, hyperscale data centres are reported to use 2.1 million litres of water a day, while smaller ones are reported to use 68,000 litres a day ([Zhang, 2024](#)).

Crucially, data centres are often located in areas with already limited water supply and exacerbate problems in those areas ([Barratt and Gambarini, 2025](#)).

Finally, there is the extraction of the raw materials needed to build servers and other data centre equipment as well as the waste they eventually become. The mining for the metals needed has its own environmental problems, and some are so-called “conflict minerals”, which means “that they are mined or traded in areas of conflict, and contribute towards perpetuating human rights abuses and armed conflict” ([Luccioni, Trevelin and Mitchell, 2024](#)).

Wang et al. (2024) predict that by 2030 generative AI could add up to 5 million metric tons of electronic waste to the global total. Given, that is a relatively small proportion of the global total, but, as experts warn, a significant one ([Crownhart, 2024](#)). Electronics often contain hazardous or toxic materials such as lead, mercury and chromium, and, if not disposed of responsibly, these can harm the environment. Another problem is the waste of valuable metals such as copper, gold and rare earth

¹⁰ It seems that parameter counts are not readily published information. The cited paper seems to be the source of the 200 billion figure floating about for GPT-4o.

elements, when electronic waste is not recycled ([Crownhart, 2024](#)). According to the 2024 Global E-Waste Monitor ([Baldé et al., 2024](#), p. 9), only 22.3% of electronic waste is formally collected and recycled in an environmentally sound manner.

I am fully aware than I am neglecting to mention the positive environmental impacts of AI, for example through helping investigating and addressing environmental problems ([UN Environment Programme, 2024](#)) as well as the steps that are being taken by authorities and cooperations to mitigate negative impacts ([Luccioni, Trevelin and Mitchell, 2024](#); [UN Environment Programme, 2024](#); [Barker, 2024](#); [Ren and Wierman, 2024](#)). Given the missing of sustainability goals by the big tech corporations ([Barker, 2025](#)) and in light of humanity's current track record in taking care of our planet, I find it not very believable that we can mitigate such a huge projected increase in resource consumption and emissions successfully. Hence, I believe it is important to showcase the massive negative environmental impact of AI clearly.

Labour exploitation and inequity

In addition to the issues regarding exploitation and inequity with mining mentioned by Luccioni, Trevelin and Mitchell ([2024](#)), there is also the problem of environmental impacts being very unfairly distributed across the planet and some regions and communities being disproportionately affected, for example, through air pollution from local fossil fuel consumption ([Ren and Wierman, 2024](#)). A particular example raised by Ren and Wierman ([2024](#)) is Google's data centre in Finland operating on 97% carbon-free energy as opposed to its ones in Asia, which only use 4%-18% carbon-free energy.

Ball ([2025](#)) also raises issues around exploitation and inequities:

"The extraction of vast amounts of data without informed consent, perpetuates a system of surveillance and control that undermines democratic principles and disproportionately affects vulnerable populations. The reliance on low-paid workers in the Global South to perform data labelling and content moderation tasks further exacerbates global inequalities, exposing these individuals to exploitative practices and precarious working conditions."

Barriers

Apart from the ethical considerations, there are also other barriers to AI use for metadata work.

AI literacy

The previously mentioned survey by Chen and Li revealed a lack of adequate training and support in relation to AI use ([Chen and Li, 2024](#), pp. 321-322), which may go some way in explaining respondents' concerns about "misunderstandings about the capabilities and limitations of AI in cataloging, which may lead to unrealistic

expectations or disappointment with the results” ([Chen and Li, 2024](#), p. 324) as well as “reservations about rushing AI integration without considering potential consequences” ([Chen and Li, 2024](#), p. 326)

This clearly shows that more needs to be done to improve AI literacy and understanding, not only for metadata professionals, but also for managers and decision makers.

Resourcing

PCC’s report on strategic planning for AI and machine learning highlights the issue of library resourcing being prohibitive for investigating, experimenting or even implementing potential AI-driven workflows:

“General concern about a lack of resources in order to investigate and implement AI. Many institutions are involved with system migrations, training for Official RDA and/or linked data, or are generally under-resourced or too small to realistically spend time working with AI.” ([Program for Cooperative Cataloguing, 2024](#), p. 2)

Few libraries around the world have the resourcing, technical expertise and equipment at their disposal to spend time on experimenting with a technology that, in order to deliver usable results, needs a deep understanding of machine learning techniques and algorithms as well as the ability to set up tools, fine-tune them to their respective needs and maintain them.

Planning, building, testing and implementing a machine learning solution such as the ones outlined earlier takes a long time, years even.

It is probably also worth saying that what works for library A does not automatically also work for library B. When it comes to metadata, we all have our local practices and idiosyncrasies to account for. Models trained on someone else’s data might not do very well with one’s own.

It would be very nice to see solutions come out of the community rather than AI-assisted cataloguing becoming yet another area where libraries need to rely on vendor-provided solutions ([Moulaison-Sandy and Coble, 2024](#), p. 381).

The Black Box

Generative AI’s “black box” nature is also a concern. Something is considered a black box when input and outputs can be seen, but how the inputs are turned into the outputs, i.e. the internal workings, are mysterious and cannot be seen ([Kosinski, no date; Bagchi, 2023](#)). Additionally, even when algorithms are known, for deep learning (which generative AI is based on), the learning process itself creates connections and patterns that mean even the creators of these processes cannot understand how they

actually work ([Kosinski, no date](#); [Metz and Weise, 2025](#); [OpenAI, 2025](#)). This means that even open-source models using deep learning are essentially black boxes.

This is problematic as it is hard to trust an output if it is not transparent how it was arrived at, and impossible to validate its path through the model. Even if the output is correct, maybe the model arrived at it for the wrong reasons (the “Clever Hans effect”). Due to the lack of understanding of the internal workings, adjusting a model that makes wrong decisions or produces bad outputs is very difficult ([Kosinski, no date](#); [Blouin, 2023](#)).

A black box model can hide security vulnerabilities. If one doesn’t know how something works, one cannot tell if it has been modified in malicious ways ([Kosinski, no date](#); [Bagchi, 2023](#)).

Black box models might also exacerbate algorithmic bias and lead to bad, maybe even outright harmful and illegal outcomes. While biases will be present in outputs if they are present in the training data, assessing if a bias exists and finding what its cause is, is especially hard in black box models ([Kosinski, no date](#); [Blouin, 2023](#)).

Lastly, Kosinski ([no date](#)) mentions trouble assessing whether one is compliant with regulations regarding the use of sensitive data in AI tools, such as for example the Artificial Intelligence Act of the European Union¹¹.

Researchers are working on improving insights into model workings, but sufficient transparency does not seem to be on the horizon ([Kosinski, no date](#))

Conclusion

All of the above may read like I oppose AI use in metadata work, but this is not the case. I am not a technophobe, and I truly believe that there are opportunities to improve record quality and to assist cataloguers and metadata managers in their work. Maybe my attitude can be best described as that of a “curator” as defined by Rosser and Hanegan ([2024](#)). In my exploration of the subject over the last couple of years I felt that, while the limitations of generative AI are mentioned and the conclusions generally align with my own here, not enough space was given to critical exploration of the technical possibilities and ethical dilemmas associated with generative AI use and this article is merely an attempt to more explicitly state the limitations and issues.

In summary, regarding the question of generative AI being able to catalogue: no, it absolutely cannot, and I believe it will not. Can machine learning catalogue? Well, maybe, but not yet. Can machine learning assist cataloguers in their work? Yes, absolutely, but the human in the loop remains a non-negotiable necessity!

¹¹ See <https://artificialintelligenceact.eu/>

I cited Moulaison-Sandy and Coble (2024) before in this article and do so again here, as they brilliantly sum up the matter:

“[N]ow is the time to look to the future and to be creative, but with a sense of the full understanding of the limitations and affordances. Yes, finding new ways in which AI can support the work of librarians, especially technical services librarians like catalogers, will be critical to future success” (Moulaison-Sandy and Coble, 2024, p. 383)

On the ethical side and the question of whether we should implement AI, I think Berkowitz (2025) has a point: we need more ethical scrutiny, policies and frameworks. For small-scale experiments and implementations, this might not be as crucial, but certainly for the adoption of AI tools, e.g. via vendor products, that rely on mainstream tools such as ChatGPT, Copilot, etc., we need to think properly about all implications and whether they outweigh the usefulness of the tool.

References

- Amram, T., Malamud, R. G., and Hollingsworth, C. (2023) Response to "From ChatGPT to CatGPT". *Information Technology and Libraries*, 42(4). Available at: <https://doi.org/10.5860/ital.v42i4.16983>
- Bagchi, S. (2023) What is a black box? A computer scientist explains what it means when the inner workings of AIs are hidden. *The Conversation*, 22 May. Available at: <https://theconversation.com/what-is-a-black-box-a-computer-scientist-explains-what-it-means-when-the-inner-workings-of-ais-are-hidden-203888> [Accessed: 24 May 2025]
- Baldé, C. P., Kuehr, R., Yamamoto, T., McDonald, R., D'Angelo, E., Althaf, S., Bel, G., Deubzer, O., Fernandez-Cubillo, E., Forti, V., Gray, V., Heart, S., Honda, S., Iattoni, G., Khetriwal, D. S., Luda di Cortemiglia, V., Lobuntsova, Y., Nnorom, I., Pralat, N. and Wagner, M. (2024) *Global E-waste Monitor 2024*. Geneva/Bonn: International Telecommunication Union and United Nations Institute for Training and Research. Available at: https://ewastemonitor.info/wp-content/uploads/2024/12/GEM_2024_EN_11_NOV-web.pdf [Accessed: 28 May 2025]
- Ball, Caroline (2025) The Unethical Underbelly of AI: A Call for Universities to Take a Stand. *UKSG News*, 586. Available at: <https://www.uksg.org/newsletter/uksg-enews-586/enews-586-editorial/> [Accessed 21 May 2025]
- Barker, C. (2024) Artificial intelligence and the environment: Taking a responsible approach. *JISC Artificial Intelligence*, 18 September. Available at: <https://nationalcentreforai.jiscinvolve.org/wp/2024/09/18/artificial-intelligence-and-the-environment-taking-a-responsible-approach/> [Accessed: 24 May 2025]
- Barker, C. (2025) Artificial intelligence and the environment: The current landscape. *JISC Artificial Intelligence*, 28 March. Available at: <https://nationalcentreforai.jiscinvolve.org/wp/2025/03/28/artificial-intelligence-and-the-environment-the-current-landscape/> [Accessed: 24 May 2025]
- Barratt, L. and Gambarini, C. (2025) Revealed: Big tech's new datacentres will take water from the world's driest areas. *The Guardian* (Online), 9 April. Available at: <https://www.theguardian.com/environment/2025/apr/09/big-tech-datacentres-water> [Accessed: 9 April 2025]

- Battersby, M. (2024) Academic authors 'shocked' after Taylor & Francis sells access to their research to Microsoft AI. *The Bookseller*, 19 July. Available at: <https://www.thebookseller.com/news/academic-authors-shocked-after-taylor--francis-sells-access-to-their-research-to-microsoft-ai> [Accessed: 24 May 2025]
- Ben Abacha, A., Yim, W.-W., Fu, Y., Sun, Z., Yetisgen, M., Xia, F. and Lin, T. (2025) *MEDEC: A Benchmark for Medical Error Detection and Correction in Clinical Notes*. Preprint. Available at: <https://doi.org/10.48550/arXiv.2412.19260>
- Berkowitz, A. E. (2025) "Slow-MO or FOMO": AI conversations at library conferences. *Public Services Quarterly*, 21(1), pp. 51-70. Available at: <https://doi.org/10.1080/15228959.2024.2442657>
- Blouin, L. (2023) AI's mysterious 'black box' problem, explained. *University of Michigan-Dearborn News*, 6 March. Available at: <https://umdearborn.edu/news/ais-mysterious-black-box-problem-explained> [Accessed: 24 May 2025]
- Brzustowicz, R. (2023) From ChatGPT to CatGPT: The Implications of Artificial Intelligence on Library Cataloging. *Information Technology and Libraries*, 42(3). Available at: <https://doi.org/10.5860/ital.v42i3.16295>
- Chen, S. and Li, M. (2024) AI for Cataloging and Metadata Creation: Perspectives and Future Opportunities from Cataloging and Metadata Professionals. *Technical Services Quarterly*, 41(4), pp. 317-332. Available at: <https://doi.org/10.1080/07317131.2024.2394919>
- Chow, E. H. C., Kao, T. J. and Li, X. (2024) An Experiment with the Use of ChatGPT for LCSH Subject Assignment on Electronic Theses and Dissertations. *Cataloging & Classification Quarterly*, 62(5), pp. 574-588. Available at: <https://doi.org/10.1080/01639374.2024.2394516>
- Chowdhury, N., Johnson, D., Huang, V., Steinhardt, J. and Schwettmann, S. (2025) Investigating truthfulness in a pre-release o3 model. Available at: <https://transluce.org/investigating-o3-truthfulness> [Accessed 24 May 2025]
- Constantino, Tor (2025) U.S. Copyright Office Shocks Big Tech With AI Fair Use Rebuke. *Forbes*, 29 May. Available at: <https://www.forbes.com/sites/torconstantino/2025/05/29/us-copyright-office-shocks-big-tech-with-ai-fair-use-rebuke/> [Accessed 31 May 2025]
- Cornish, B. and Scott, B. (2025) Identify, Obtain, Explore: using NLP to link article and journal records in the NHM library catalogue. *Catalogue & Index*, 211, pp. 20-28. Available at: <https://journals.cilip.org.uk/catalogue-and-index/article/view/749>
- Corrado, E. M. (2021) Artificial Intelligence: The Possibilities for Metadata Creation. *Technical Services Quarterly*, 38(4), pp. 395-405. Available at: <https://doi.org/10.1080/07317131.2021.1973797>
- Crawford, K. (2024) Generative AI's environmental costs are soaring - and mostly secret. *Nature*, 626, p. 693. Available at: <https://doi.org/10.1038/d41586-024-00478-x>
- Crownhart, C. (2024) AI will add to the e-waste problem. Here's what we can do about it. *MIT Technology Review*, 28 October. Available at: <https://www.technologyreview.com/2024/10/28/1106316/ai-e-waste/> [Accessed: 24 May 2025]
- Deutsche Nationalbibliothek (2023) *Launch of cataloguing machine EMa*. Available at: https://jahresbericht.dnb.de/Webs/jahresbericht/EN/2022/Hoehepunkte/Erschliessungsmaschine/erschliessungsmaschine_node.html [Accessed: 24 May 2025]
- DeZelar-Tiedman, C. (2023) Response to "From ChatGPT to CatGPT". *Information Technology and Libraries*, 42(4). Available at: <https://doi.org/10.5860/ital.v42i4.16991>
-

- Eaton, L. (2024) Research Insights #12: Copyrights and Academia: Scholarly authors are not going to be happy... *AI+Edu=Simplified*, 23 July. Available at: <https://aiedusimplified.substack.com/p/research-insights-12-copyrights-and> [Accessed: 24 May 2025]
- ExLibris (no date a) *AI Bibliographic Records Enrichment*. Available at: https://knowledge.exlibrisgroup.com/Content/Knowledge_Articles/Alma/Knowledge_Articles/AI_Bibliographic_Records_Enrichment [Accessed: 31 May 2025]
- ExLibris (no date b) *AI Metadata Enrichment for Libraries*. Available at: https://knowledge.exlibrisgroup.com/Alma/Product_Materials/010Roadmap/AI_Metadata_Enrichment_for_Libraries [Accessed: 31 May 2025]
- Floyd, D. (2023) Response to "From ChatGPT to CatGPT". *Information Technology and Libraries*, 42(4). Available at: <https://doi.org/10.5860/ital.v42i4.16995>
- Gallegos, I. O., Rossi, R. A., Barrow, J., Tanjim, M. M., Kim, S., Deroncourt, F., Yu, T., Zhang, R. and Ahmed, N. K. (2024) Bias and Fairness in Large Language Models: A Survey. *Computational Linguistics*, 50(3), pp. 1097-1179. Available at: https://doi.org/10.1162/coli_a_00524
- Goldman Sachs (2024) *AI, data centers and the coming US power demand surge*. Available at: <https://www.goldmansachs.com/pdfs/insights/pages/generational-growth-ai-data-centers-and-the-coming-us-power-surge/report.pdf> [Accessed: 24 May 2025]
- Golub, K., Suominen, O., Mohammed, A.T., Aagaard, H. and Osterman, O. (2024) Automated Dewey Decimal Classification of Swedish library metadata using Annif software. *Journal of Documentation*, 80(5), pp. 1057-1079. Available at: <https://doi.org/10.1108/JD-01-2022-0026>
- Google (no date) *What is Retrieval-Augmented Generation (RAG)?* Available at: <https://cloud.google.com/use-cases/retrieval-augmented-generation> [Accessed: 03 June 2025]
- Hicks, M. T., Humphries, J. and Slater, J. (2024) ChatGPT is bullshit. *Ethics and Information Technology*, 26(2). Available at: <https://doi.org/10.1007/s10676-024-09775-5>
- Inkinen, J., Lehtinen, M. and Suominen, O. (2025) *Annif Users Survey: Understanding Usage and Challenges*. National Library of Finland. Available at: <https://urn.fi/URN:ISBN:978-952-84-1301-1> [Accessed: 24 May 2025]
- Kosinski, M. (no date) *What is black box artificial intelligence (AI)?* Available at: <https://www.ibm.com/think/topics/black-box-ai> [Accessed: 31 May 2025]
- Library of Congress (no date) *Exploring Computational Description*. Available at: <https://labs.loc.gov/work/experiments/ECD> [Accessed: 24 May 2025]
- Lowagie, H. (2024) Harnessing Power Apps and AI for Automated Cataloguing: Innovations in Bibliographic Record Creation. *Catalogue & Index*, 209. Available at: <https://journals.cilip.org.uk/catalogue-and-index/article/view/697> [Accessed: 01 May 2025]
- Luccioni, A. S., Jernite, Y. and Strubell, E. (2024) Power Hungry Processing: Watts driving the cost of AI deployment?, *FACt '24: The 2024 ACM Conference on Fairness, Accountability, and Transparency*, Rio de Janeiro, 3-6 June. Available at: <https://doi.org/10.1145/3630106.3658542>
- Luccioni, A. S., Viguier, S. and Ligozat, A. L. (2023) Estimating the carbon footprint of BLOOM, a 176B parameter language model. *Journal of Machine Learning Research*, 24(253), pp. 1-15. Available at: <https://jmlr.org/papers/v24/23-0069.html> [Accessed: 31 May 2025]
- Luccioni, S., Trevelin, B. and Mitchell, M. (2024) The Environmental Impacts of AI - Policy Primer. *Hugging Face Blog*, 3 September. Available at: <https://doi.org/10.57967/hf/3004>

- Metz, C. and Weise, K. (2025) AI is getting more powerful, but its hallucinations are getting worse. *New York Times* (Online), 5 May. Available at: <https://www.nytimes.com/2025/05/05/technology/ai-hallucinations-chatgpt-google.html> [Accessed: 31 May 2025]
- Mollick, E. (2025 a) The End of Search, The Beginning of Research: The first narrow agents are here. *One Useful Thing*, 3 February. Available at: <https://www.oneusefulthing.org/p/the-end-of-search-the-beginning-of> [Accessed 6 April 2025]
- Mollick, E. (2025 b) A new generation of AIs: Claude 3.7 and Grok 3. *One Useful Thing*, 24 February. Available at: <https://www.oneusefulthing.org/p/a-new-generation-of-ais-claude-37> [Accessed: 6 April 2025]
- Moulaison-Sandy, H. and Coble, Z. (2024) Leveraging AI in Cataloging: What Works, and Why? *Technical Services Quarterly*, 41(4), pp. 375-383. Available at: <https://doi.org/10.1080/07317131.2024.2394912>
- Murgia, M., Clark, D., Learner, S., de la Torre Arenas, I., Joiner, S., Hemingway, E. and Hawkins, O. (2023) Generative AI exists because of the transformer. *Financial Times*, 12 September. Available at: <https://ig.ft.com/generative-ai/> [Accessed: 24 May 2025]
- National Library of Finland (2025) *Annif*. Available at: <https://annif.org/> [Accessed: 24 May 2025]
- OCLC (2025) *Implementing AI to further scale and accelerate WorldCat de-duplication*. Available at: <https://www.oclc.org/en/news/announcements/2025/ai-worldcat-deduplication.html> [Accessed: 24 May 2025]
- OpenAI (2025) *OpenAI o3 and o4-mini System Card*. Available at: <https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf> [Accessed: 1 June 2025]
- Poley, C., Uhlmann, S., Busse, F., Jacobs, J.-H., Kähler, M., Nagelschmidt, M. and Schumacher, M. (2025) Automatic Subject Cataloguing at the German National Library. *Liber Quarterly*, 35. Available at: <https://doi.org/10.5337/lq.19422>
- Program for Cooperative Cataloguing (2024) *PCC Task Group on Strategic Planning for AI and Machine Learning: Final Report Transmittal & Tracking Sheet*. Available at: <https://www.loc.gov/aba/pcc/taskgroup/TG-Strategic-Planning-AI-final-report.pdf> [Accessed: 21 May 2025]
- Resnik, P. (2025) Large Language Models Are Biased Because They Are Large Language Models. To be published in *Computational Linguistics* [Peer-reviewed accepted version]. Available at: https://doi.org/10.1162/coli_a_00558
- Reisner, A. (2025) The Unbelievable Scale of AI's Pirated-Books Problem. *The Atlantic* (Online), 20 March. Available at: <https://www.theatlantic.com/technology/archive/2025/03/libgen-meta-openai/682093/> [Accessed: 28 May 2025]
- Ren, S. and Wierman, A. (2024) The Uneven Distribution of AI's Environmental Impacts. *Harvard Business Review*, 15 July. Available at: <https://hbr.org/2024/07/the-uneven-distribution-of-ais-environmental-impacts> [Accessed: 24 May 2025]
- Rosser, C. and Hanegan, M. (2024) Cyborgs and Centaurs, Prophets and Priests: Anywhere Left for Curators? *Atla*, 17 April. Available at: <https://www.atla.com/blog/cyborgs-and-centaurs-prophets-and-priests-anywhere-left-for-curators/> [Accessed: 28 May 2025]
- Saccucci, C. and Potter, A. (2024 a) Exploring Computational Description: LC Labs Planning Framework in Action, *OCLC RLP Webinar*, 12 March. Available at: <https://www.oclc.org/research/events/2024/ai-planning-framework-in-action.html> [Accessed: 15 May 2025]

- Saccucci, C. and Potter, A. (2024 b) Exploring Machine Learning: A Cataloging Experiment at the Library of Congress, *PCC Operations Committee Meeting*, 2 May. Available at: <https://www.loc.gov/aba/pcc/documents/OpCo-2024/Potter-Saccucci-Machine-Learning.pdf> [Accessed: 24 May 2025]
- Shojaee, P., Mirzadeh, I., Alizadeh, K., Horton, M., Benigo, S. and Farajtabar, M. (2025) *The Illusion of Thinking: Understanding the Strengths and Limitations of Reasoning Models via the Lens of Problem Complexity*. Available at: <https://ml-site.cdn-apple.com/papers/the-illusion-of-thinking.pdf> [Accessed: 10 June 2025]
- Strubell, E., Ganesh, A., and McCallum, A. (2020) Energy and policy considerations for modern deep learning research. *Proceedings of the AAAI conference on artificial intelligence*, 34(9), pp. 13693-13696. Available at: <https://doi.org/10.1609/aaai.v34i09.7123>
- Suominen, O., Inkinen, J., Virolainen, T., Fürneisen, M., Kinoshita, B. P., Veldhoen, S., Sjöberg, M., Zumstein, P., Neatherway, R. and Lehtinen, M. (2023) *Annif* (v1.0.0). Zenodo. <https://doi.org/10.5281/zenodo.8262313>
- Taniguchi, S. (2024) Creating and Evaluating MARC 21 Bibliographic Records Using ChatGPT. *Cataloging & Classification Quarterly*, 62(5), pp. 527-546. Available at: <https://doi.org/10.1080/01639374.2024.2394513>
- Uhlmann, S. and Grote, C. (2021) Automatic subject indexing with Annif at the German National Library (DNB). *Semantic Web in Libraries*, Online, 29 November – 3 December. Available at: <https://swib.org/swib21/slides/03-02-uhlmann.pdf> [Accessed: 24 May 2025]
- UN Environment Programme (2024) *AI has an environmental problem. Here's what the world can do about that*. Available at: <https://www.unep.org/news-and-stories/story/ai-has-environmental-problem-heres-what-world-can-do-about> [Accessed: 24 May 2025]
- United States Copyright Office (2025) *Copyright and artificial intelligence Part 3: Generative AI Training*. Pre-publication version. Washington: United States Copyright Office. Available at: <https://www.copyright.gov/ai/Copyright-and-Artificial-Intelligence-Part-3-Generative-AI-Training-Report-Pre-Publication-Version.pdf> [Accessed: 1 June 2025]
- Wang, P., Zhang, L.-Y., Tzachor, A. and Chen, W.-Q. (2024) E-waste challenges of generative artificial Intelligence. *Nature Computational Science*, 4, pp. 818-823. Available at: <https://doi.org/10.1038/s43588-024-00712-6>
- Weinryb-Grohsgal, L., Potter, A. and Saccucci, C. (2024) Could Artificial Intelligence Help Catalog Thousands of Digital Library Books? An Interview with Abigail Potter and Caroline Saccucci. *The Signal*, 19 November. Available at: <https://blogs.loc.gov/thesignal/2024/11/could-artificial-intelligence-help-catalog-thousands-of-digital-library-books-an-interview-with-abigail-potter-and-caroline-saccucci/> [Accessed: 24 May 2025]
- Wu, J., Gan, W., Chen, Z., Wan, S. and Yu, P. S. (2023) Multimodal Large Language Models: A Survey. *2023 IEEE International Conference on Big Data*, Sorrento, 15 – 18 December. Available at: <https://doi.org/10.1109/BigData59044.2023.10386743>
- York, E. and Hanegbi, D. (2024) *Metadata Enrichment using AI First Glance at Research and Findings*. Available at: https://knowledge.exlibrisgroup.com/@api/deki/files/166843/AI_enrichment_-_March_27.pdf?revision=1 [Accessed: 24 May 2025]
- Zewe, A. (2025) Explained: Generative AI's environmental impact. *MIT News*, 17 January. Available at: <https://news.mit.edu/2025/explained-generative-ai-environmental-impact-0117> [Accessed: 24 May 2025]


Zhang, M. (2024) *Data Center Water Usage: A Comprehensive Guide*. Available at: <https://dgtlinfra.com/data-center-water-usage/> [Accessed 24 May 2025]

Book review: The DEI Metadata Handbook

Reviewed by: **Ashleigh Weir**

Library Administrator – Metadata Services, Sheffield Hallam University

Received: 02 June 2025 | Published: 17 June 2025

Wintermute, H.E., Campbell, H.M., Dieckman, C.S., Rose, N.L. and Thulsidhos, H. (2024) *The DEI Metadata Handbook: A Guide to Diverse, Equitable, and Inclusive Description*. Ames, Iowa: Iowa State University Digital Press. ISBN 978-1-958291-09-2 (online), ISBN: 978-1-958291-10-8 (print), DOI [10.31274/isudp.2024.153](https://doi.org/10.31274/isudp.2024.153) 

CONTACT Ashleigh Weir  A.A.Weir@shu.ac.uk  Sheffield Hallam University

The DEI Metadata Handbook is a well-considered and practical guide to approaching issues of diversity, equality and inclusion in metadata. The book was authored by Harriet E. Wintermute, Heather M. Campbell, Christopher S. Dieckman, Nausicaa L. Rose and Hema Thulsidhos, all metadata librarians at Iowa State University at time of writing, with Harriet E. Wintermute now working as the Chair of Acquisitions, Cataloging, Metadata, and E-Resources at the University of Nebraska–Lincoln Libraries. Whilst this book is aimed mostly at library metadata creators, it could also apply to any information professionals who may create resource descriptions in their work, and also be of interest to library and archive students. Because of this, the book feels like a great introduction to the subject for someone new to EDI in metadata. The handbook aims to help readers gain awareness of EDI related issues and learn new techniques to review both existing metadata and improve metadata for diverse ranges of resources going forward.

A point put across in this book multiple times that I really valued was that, while the initialism of DEI in metadata is relatively new, the goals of it, or in libraries in general, are grounded in historical origins. Looking at DEI in metadata aims to “enhanc[e] diverse representation in descriptive metadata; improv[e] discovery of diverse resources; and mitigat[e] negative effects of inaccurate, outdated, or offensive terminology” (p. 1). Alongside contemporary practice, the handbook gives multiple examples of historical approaches to presenting a diverse range of resources in the library from up to a hundred years ago. Knowing that this work has always been important for some cataloguers really gives it a historical grounding, which counters any ideas that EDI is just ‘another trend’ in the library world.

The first chapter covers inclusive description in free-text fields when cataloguing. It prompts the metadata creator to think about how and when you should refer to identity, looking at relevancy, users’ needs, accuracy and respect, four considerations

which are listed multiple times through the book. The chapter helpfully looks at different aspects of identity such as ability, class, gender and sexuality, race and ethnicity, and religion. It gives useful and up-to-date pointers on what language may be harmful or outdated for different communities and what language to use instead. One main piece of advice is to not assume someone's identity based on images or names, this falling under their point about accuracy; for example, to not list someone's race or gender unless you are sure that is an identity they publicly hold. This chapter also gives examples of how different library collections have used content warnings in the form of harmful language statements to alert users of outdated or harmful language that may be in the collection, either in older library material or the library records themselves.

The second chapter looks at name authorities and how to use them ethically. It gives examples of where authors have used pseudonyms and anonymously authored texts. It also discusses married women's names and transgender creators where they do not wish to have their deadname (previous name) known publicly. A main point it adds is, where the author is still living, if there is uncertainty about a name, to try to contact them if possible. This chapter also gives useful examples of how to transcribe names written in non-Latin script.

The third chapter of the book focuses on descriptions of people and groups, including creators and contributors, audiences, depicted people, and others. As in previous chapters, it speaks about accuracy but also respect and privacy; it explains the importance of respecting a person's request that information about them be removed from a library record and how to approach that. This chapter also discusses Library of Congress Subject Headings (LCSH) being problematic for not accounting for all possible groups. For example, it has many instances of markedness, where linguistically it marks something as other with language, like 'doctors' versus 'women doctors'. The authors of the handbook suggest a range of approaches to counter this, such as making full use of the Library of Congress Demographic Group Terms (LCDGT) and employing other controlled vocabularies such as Homosaurus, FAST and ERIC. The chapter gives examples of how one could use these. Whilst a lot of the book is based on MARC, it also gives examples of how to do so in BIBFRAME and the CIDOC Conceptual Reference Model.

Chapter four discusses classification and the biases that can be embedded within classification systems such as the Library of Congress Classification (LCC) or Dewey Decimal Classification. I really valued how the authors look at both working within DDC and LCC and also working with alternative classification systems. The chapter gives examples of how librarians in the past have created new classification systems that work better for their collection and their users, and also gives practical questions to ask yourself when looking into implementing a new classification scheme. This could be especially useful for librarians with small or specialist collections.

The fifth chapter explores DEI with regards to subject headings and how to improve the inclusivity of subject metadata. Similar to chapter three, it explores how to work within LCSH but also additional vocabularies you could use and how to implement them with helpful questions to ask yourself. It gives a helpful starting point for how to go about proposing a new LCSH heading through the Subject Authority Cooperative Program (SACO). It gives a timeline and an example of the high-profile case spanning 2013-2021 regarding the subject heading 'illegal aliens' and eventually getting this replaced. This chapter also includes a case study from the authors of the book, at Iowa State University itself, where they worked with Iowa Indigenous peoples to create vocabularies for Indigenous people that use terms used by the communities.

The last chapter looks at accessibility as it relates to metadata. It points out the importance of making note of accessibility features of physical library materials, despite it often being more common for digital materials. This would include adding accessibility content (341 in MARC) and accessibility notes (532) such as information about captions, audio descriptions, braille and more. The chapter also gives a substantial list of best practices and dos and don'ts for non-physical resources including alt text, extended descriptions, audio descriptions, captions, and transcripts for both MARC and Dublin Core. The chapter concludes with a powerful remark about the importance of making use of these features: "Besides improving accessibility, a meticulous approach to metadata serves as a powerful tool to promote inclusivity and foster an equitable and just information environment for diverse communities" (p.84).

The book opens with a quote from the Cataloguing Ethics Steering Committee's Cataloguing Code of Ethics that establishes the crucial role of cataloguers in deciding how information is represented, and ultimately I think that this book does well to show how we can all make considerations in our metadata work to improve representation and discovery of a diverse range of materials for our library users. Whilst a lot of the content in here may not be brand new information for anyone with some awareness of EDI in metadata, the handbook feels like a useful guide in that it has synthesised a lot of knowledge on the topic into a digestible, quick-to-read book with excellent starting points for anyone looking into making these changes in their library. Each chapter is full of citations, practical examples, case studies and concludes with a list of resources for further reading, and for this reason I think it would be a great addition to any metadata or cataloguing team's set of tools.



Book review: RDA and Serials Cataloging

Reviewed by: **Natasha Aburrow-Jones**

Systems Librarian, Advocates Library

Received: 06 June 2025 | Published: 17 June 2025

Jones, E. (2025) *RDA and Serials Cataloging*. Second edition. Chicago: ALA Editions. ISBN 978-0-8389-4871-2 (print), ISBN 979-8-89255-539-5 (PDF), ISBN 979-8-89255-538-8 (ePub)

CONTACT Natasha Aburrow-Jones  tashaaburrowjones@gmail.com  Advocates Library

This book by Ed Jones is your key to demystifying that trifecta of complicated cataloguing scenarios – Official RDA, serials, and MARC21. An updated version of his original book, *RDA and Serials Cataloging*, first published in 2013, this second edition covers all the exciting developments in the cataloguing world over the past decade.

Serials are fascinating items to describe – they change and shift over the course of time, but still maintain their identity. They're seen as "difficult" to describe, and are treated with trepidation. I have been involved in serials cataloguing for over thirty years, and I found myself relaxing into this book, knowing that it was easy to read, full of practical advice, and not afraid to shy away from complicated scenarios and examples. As the author says:

"RDA and Serials Cataloging is designed to be used by serials catalogers who are new to RDA and by monograph catalogers who are new to serials cataloging"

The first part is split into three chapters, tackling an overview of serials, an introduction to RDA, and searching serials. There was a lot on the history of serials and their description, which I found absolutely riveting, starting off with that question asked of all those involved in serials – "What is a serial?" There is no definitive answer (as I have heard before, everything is a serial), but:

"... a serial is like a good work of art: you may not be able to say what a serial is, but you know one when you see one."

Setting the ground work of how serials behave is a good introduction to their description, and this book covers print, electronic, continuing and integrating resources, both scholarly and popular, and everything in between. The history of serials cataloguing is something that is worth a read, and Ed Jones gives a succinct description, and then treats the development of RDA in the same way.

The second part of the book is entitled "*Cataloging Serials and Ongoing Integrating Resources Using RDA*". And it really is very thorough! Most of this part is about cataloguing in MARC21, the carrier standard still used by the vast majority of libraries. So, we have a combination of how to treat a serial in both RDA and MARC21, covering content and carrier. Examples are plentiful and comprehensive, and so many variables are covered. If you don't have access to the RDA Toolkit, this book is a good manual to use in its place (with the caveat that the RDA Toolkit will be constantly updated, and this defined book will not – well, not until the next edition, at any rate). I liked the approach of thinking about how to tackle a serial before you start filling in your serials bibliographic record – there is discussion around major and minor title changes, a breakdown of the 007 and 008 fields in MARC21 (Physical Description Fixed Field, and Fixed Length Data Elements, respectively), when to create a new record... all the points that can look rather intimidating to anyone new to serials cataloguing. There are some useful points throughout marking the difference between AACR2, ISBD and RDA, which are useful for those who are coming new to RDA, but have catalogued in AACR2. For example (about editions):

"Prescribed punctuation:

AACR2: Edition area ends with a single period /stop.

ISBD (revised): Edition area ends with two periods /stops whenever the last element in the area ends in a period /stop (e.g., an abbreviation).

RDA: Does not prescribe punctuation."

"*RDA and Serials Cataloging*" starts off gently with print serials, but then does move into the heavy hitting world of online serials, born digital, version control, integrating resources, different formats (print and electronic, and even microform), and shows with full record examples ways of dealing with these. Reading the theory is one thing, but backing it up with solid examples helps massively when it comes to cataloguing serials, or, as RDA calls them, diachronic works.

There is also discussion about carrier standards other than MARC21, such as BIBFRAME. RDA was always meant to be a standard used more widely than purely in the library world, and other forms of data description will affect how we describe items in other ways. Linked Data, mapping between different vocabularies and standards – all of these are referenced as something to explore and look forward to. As the author says:

[The future] "holds both great promise and novel challenges, and libraries can expect to be in the thick of it. A brave new world indeed."

To conclude, this volume should be a valued addition to any cataloguing bookshelf. It describes how to handle serials cataloguing in a way that is easy to read and makes sense. Serials should not be shrouded in mystery, and "*RDA and Serials Cataloging*" is a very useful tool to combat any anxiety associated with serials cataloguing.

Book review: Many Pathways for Discovery

Reviewed by: **James Clark**

Southampton Solent University

Received: 09 June 2025 | Published: 17 June 2025

Mullin, Casey A. (2024) *Many Pathways for Discovery: Describing Music Resources Using Faceted Vocabularies*. Middleton, WI: MLA and A-R Editions. ISBN 978-0-89579-911-1 (print), ISBN 978-0-89579-912-8 (online)

CONTACT James Clark  james.clark@solent.ac.uk  Southampton Solent Univeristy

An important aspect of cataloguing notated or performed music resources is that subject analysis requires a description of what the resource is as opposed to what it is about; cataloguers must record the genre or form of the work(s), the instruments used or required and their number. Such information can usually be expressed via LCSH, yet music cataloguers have for many years sought to record this information separately in a more machine-readable way: since at least the 1980s the coded fields 045, 047, and 048 have been used by some libraries¹. This effort has been given renewed impetus by the development, begun in 2009, of more comprehensive thesauri for music: the *Library of Congress Medium of Performance Thesaurus for Music* (LCMPT) and musical elements of the *Library of Congress Genre/Form Terms for Library and Archival Materials* (LCGFT). In this context, there is a clear need for a cataloguing manual to advise on the application of these fields. Casey Mullin is well-qualified to write such a manual – he was involved in the development of LCMPT and LCGFT – and with this slim volume he fills the need well.

As well as explaining how to describe medium of performance (LCMPT in MARC 382) and genre/form (LCGFT in 655 and 380), Mullin also explains how to record chronological information (046 and 388), geographic information (370) and demographic characteristics of audience and contributors (Library of Congress Demographic Group Terms in 385 and 386). This approach is more comprehensive and allows for describing greater complexity than the encoded data that was once used (e.g., in 008, 045, 047, and 048), and it renders such data redundant. But the new vocabularies are also intended as a means “to overcome the limitations and drawbacks of LCSH” (p. 26), and their developers “have envisioned a future where LCSH practices ... that are duplicative of these faceted terms will eventually be phased out” (p. 27), though for the time being, best practice is for cataloguers to use both the new terms and LCSH.

¹ The Library of Congress ceased using these fields in 1991, but their use persisted in some libraries (Holden et al., 2019, pp. 597-598).

The main body of the book, chapters 2 to 5, serves as a manual for the application of these vocabularies. A reader not already familiar with music cataloguing will need to use this book in conjunction with more general works (p.4), and all readers will need to consult the thesauri directly and have knowledge of existing guidance², but this book should help both experienced and new music cataloguers begin using the vocabularies. The writing is clear and concise, both when issuing authoritative guidance and when acknowledging the need for cataloguer judgement, e.g., “When performers play multiple instruments, use 382 subfield \$d to indicate this doubling relationship. Input a numeral in subfield \$n after each subfield \$a and \$d” (p. 34), and “Use judgement in determining how many instruments are feasible to record” (p. 35). His advice (on p. 44) about determining how to code the indicators in a particularly tricky case might benefit cataloguers in many other contexts; it is simply: “Do not agonize.”

As you would expect from a cataloguing manual, there is barely a page that does not include at least one example illustrating the points made. As well as the field under discussion, the examples typically include the 245 and 650 fields. Many sections have a subheading *LSCH comparison*, which serves to guide the user in adding LCSH alongside faceted fields; cumulatively, these short paragraphs also illustrate the limitations of LCSH as compared to the new vocabularies. Another recurring heading is the *Retrospective implementation note*. This section references what “an automated program” (*passim*) might generate and highlights what further action cataloguers would need to take to remediate the automated work. Despite the generic phrasing, these sections are surely written with the OCLC Music Toolkit in mind. Cataloguers working outside OCLC, and within technological limitations, might have to work with a retrospective tool that produces very different results, and thus they would not be able to use these sections as they are intended. However, their inclusion does at least emphasise to cataloguers the possibility of retrospective conversion, and if it prompts someone to press their systems team for a program, then it will have served a purpose.

Another useful recurring feature is the ‘extracts’ from the thesauri, such as the listing of instrument category terms (pp. 37-38). What Mullin has done here is extract all the category terms (i.e., not terms for specific instruments) for instruments from the thesaurus and arrange them hierarchically (as opposed to the alphabetic arrangement in LCMPT), e.g., aerophone > wind instrument > brass instrument. Therefore this ‘extract’ does not replicate what a cataloguer would find in the thesaurus, but presents a new arrangement of the thesaurus, one that is useful for the specific context and that aids a deeper understanding of the thesaurus’ structure.

Chapters 2-4 discuss the application of the new vocabularies by the type of music being described, respectively: instrumental music (in the western art music tradition); vocal music (in the western art music tradition); popular and folk/traditional music.

² See, for example (all freely available online), [Library of Congress \(2022\)](#), [Music Library Association \(2023\)](#), and [Music Library Association \(2024\)](#).

Duplication of content is reduced by references to earlier chapters. This is sensible, but it might inconvenience someone using the book as a reference resource. Chapter 5 explains how to describe the content type of musical resources, and so naturally addresses them according to this feature, with sections on scores, audio recordings, and video recordings.

In the final two chapters, co-written by Kevin Kishimoto, the book departs from its cataloguing manual format and offers instead brief introductions to related topics. Chapter 6 introduces non-MARC metadata and describes how the vocabularies might be applied in BIBFRAME. Chapter 7 describes how the new vocabularies might be leveraged in a discovery environment and includes some suggestions that can be applied now (in some environments) and some that should be possible in the future. As in the earlier chapters, examples abound, but I wonder if they are so useful here. Does the reader really need three examples, about two and a half pages of an 11-page chapter, of Alma discovery normalisation rules when they are all freely available online? Similarly, Chapter 6 includes 18 examples of BIBFRAME data in N-Triple or Turtle serialisation. A cataloguer working in BIBFRAME need not encounter these formats, and while I can see how some exposure to them could aid understanding, for the examples to comprise almost a third of the chapter seems excessive. This is a minor quibble, but such use of space is particularly noticeable in a book whose main text runs to just 141 pages. However, it is certainly important that cataloguers have at least a basic understanding of the topics covered in these two chapters, and Mullin and Kishimoto do a sound job of giving a brief – self-consciously so (pp. 119, 131) – introduction to the topics.

The book is generally well produced. The binding seems secure, but it is loose enough to allow the pages to remain open when the book is put down, a feature not to be under-estimated in a reference work! But the book would have been easier to use for reference if the running headers had indicated what vocabulary was under discussion on each page. I think there are misprints in the examples on pages 61 and 123, but neither would cause confusion.

I came to this book with an interest in faceted description and discovery, but no knowledge of the vocabularies developed for music resources and little experience of cataloguing music. It is difficult for me to say, therefore, whether cataloguers who have already been using the new vocabularies (and perhaps those working in libraries that have prepared substantial internal documentation on them) would benefit from this book. But no doubt the primary intended audience for this book is cataloguers coming to the topic with little or no prior understanding of it, and I would not hesitate to recommend it to that group.

References

Holden, C., Knop K. and Newcomer, N. (2019) 'Music Discovery: Past, Present, and Future,' *Notes*, 75(4), pp. 591-619. Available at: <https://doi.org/10.1353/not.2019.0007>

- Library of Congress (2022) *Medium of Performance Thesaurus for Music Manual*. Available at: <https://www.loc.gov/aba/publications/FreeLCMPT/freelcmpt.html> [Accessed: 9 June 2025]
- Music Library Association (2023) *Best Practices for Using LCMPT*. Available at: <https://cmc.wp.musiclibraryassoc.org/documents/best-practices-for-using-lcmpt/> [Accessed: 9 June 2025]
- Music Library Association (2024) *Best Practices for Using LCGFT for Music Resources*. Available at: <https://cmc.wp.musiclibraryassoc.org/documents/best-practices-for-using-lcgft-for-music-resources/> [Accessed: 9 June 2025]



Catalogue & Index is electronically published by the Metadata & Discovery Group of the Chartered Institute of Library and Information Professionals (CILIP) (Charity No. 313014).

Submissions: Please follow the submission guidance on our website.

Book reviews: Please contact the editors.

Advertising rates: GBP 70.00 full-page; GBP 40.00 half-page. Prices quoted without VAT.

Editors: Karen F. Pierce and Fran Frenzel

ISSN 2399-9667

MDG Website: www.cilip.org.uk/mdg

MDG blog: <http://catandindexgroup.wordpress.com/>

