

Who wrote this??

And who is an author, anyway?

AI authors in the catalog

Sarah Hovde
UMD Libraries Research and Innovative Practice Forum
June 5, 2024

Hello! I am here to talk about AI in the library.

But not in the information literacy sense, or the just experimenting with it sense, or in the "it's going to put us all out of work!" sense. I am here to talk about how catalogers can describe resources created by an AI, and more specifically, whether AIs can be treated as "authors."

When I say AI, which "AI" am I talking about? My focus here is on generative AIs that are used to create texts, images, video, music, and probably even more, and in particular at text created by large language models, such as ChatGPT and Gemini (formerly Bard), which have seen a vast public surge of interest in the last few years. The history of generative text programs goes back somewhat further, though: story-telling computer programs existed as early as the 1960s and 1970s*; digital poetry programs were created in the 1960s in languages such as BASIC, TRAC, APL, and FORTRAN**; and the first book widely credited to be "written by an AI" was published in 1984.*

*Leah Henrickson. "The policeman's beard is algorithmically constructed," *3:am Magazine*, July 16, 2018, viewed June 1, 2024.

<https://www.3ammagazine.com/3am/the-policemans-beard-is-algorithmically-constructed/>

**C.T. Funkhouser. "Chapter 1: Origination: Text Generation" in *Prehistoric digital poetry: an archaeology of forms, 1959-1995* (Tuscaloosa: University of Alabama Press, 2007), 31-84.

AI in our catalog?

- Exoanthropology : ꞑ Dialogues with AI / ꞑ Robert Leib.
- Voidopolis / ꞑ Kat Mustatea ; afterwords by Charlotte Kent and Arielle Saiber.
 - created with a "modified GPT-2 text generator"
- Benny the blue whale : ꞑ a descent into story, language and the madness of ChatGPT / ꞑ Andy Stanton vs ChatGPT.
- Regular expression puzzles and AI coding assistants : ꞑ 24 puzzles solved by the author, with and without assistance from Copilot, ChatGPT and more / ꞑ David Q. Mertz.
- Melt / ꞑ Philip Zimmermann.
 - "a conversation with two different generative artificial intelligence entities. ChatGPT and DALL-E, both from OpenAI: one generates the text, the other the pictures."
- The Policeman's beard is half-constructed : computer prose and poetry / by Racter ; illustrations by Joan Hall ; introduction by William Chamberlain.
- Divine Transplant: "Uprooted for Growth"
 - 7101_ ꞑ Midjourney AI V5.2 2023, ꞑ cover design by. ꞑ oth

In our present day of LLMs, several studies and news reports have suggested that there is a steadily growing amount of books and articles that are created partially or fully with AI. It was hard to find an estimate of how widespread this phenomenon is - a Reuters article from last year noted about 200 books on Amazon with ChatGPT listed as an author or co-author as of February 2023.* However, that is not counting books and articles that are written with uncredited AI but passed off as the work of humans alone. A recent study by librarian Andrew Gray estimated that as many as 60,000 papers published in 2023 may have been created with assistance from an LLM, based on sharp increases in the prevalence of certain keywords often associated with AI-generated text.**

A quick and non-comprehensive search of Alma shows about half a dozen titles whose creation involved some kind of AI process, plus more records in the Community Zone (which includes all Alma libraries). These are the resources that **self-identify** as having been created by or with AI, and that I was able to find in about five minutes of searching.

*Greg Bensinger. "Focus: ChatGPT launches boom in AI-written e-books on Amazon," *Reuters*, February 21, 2023, viewed June 1, 2024.

<https://www.reuters.com/technology/chatgpt-launches-boom-ai-written-e-books-amazon-2023-02-21/>

**Andrew Gray, "ChatGPT "contamination": estimating the prevalence of LLMs in the scholarly literature," *arXiv*, March 25, 2024, viewed June 1, 2024.

<https://doi.org/10.48550/arXiv.2403.16887>

Bibliographic relationships

- relationships to other resources
 - translation, sequel, parody adaptation
- relationships to persons, corporate bodies, etc.
 - author, publisher, illustrator, cinematographer
- topical relationships
 - subject, genre

How are these materials represented in a catalog? As a quick refresher: catalogers create bibliographic records, which include bits of metadata such as the title or the number of pages which describe the resource itself. Other pieces of metadata in the record explain the resource's relationships: how it is connected to other resources; its connections to people, families, corporate bodies involved in the resource's creation or publication; or topical relationships such as its subject and genre. These relations are indicated through a combination of MARC field coding, controlled vocabulary terms, and sometimes free text notes. So for the most part, a text generated by an LLM can be treated as any other text in the bibliographic record, with its title, number of pages, and other bits of metadata added accordingly. It's the relationships, and specifically the author relationship, that gets tricky.

200 years of cataloging standards

- Anthony Panizzi (British Library), 91 Rules for Compilation of the Catalogue, 1841
- Charles Ammi Cutter (Boston Athenaeum, Rules for a dictionary catalogue, 1876
- American Library Association, Catalog Rules, Author and Title Entries, 1908
- ALA, A.L.A. Catalog rules: author and title entries (revision/preliminary 2nd edition of 1908 edition), 1941
- ALA, A.L.A. cataloging rules for author and title entries, 2nd edition, 1949
- American Library Association, Library of Congress, Library Association, and Canadian Library Association, Anglo-American Cataloging Rules (AACR), 1967
- Anglo-American Cataloging Rules, 2nd edition (AACR2), 1978
- RDA Steering Committee. Original RDA Toolkit, 2013
- RDA Steering Committee. Official RDA Toolkit, 2020

The author relationship is usually recorded in a MARC 100 or 110 field, sometimes in a 700 or 710 field, and it is a major way that users look for works in the catalog. Consequently, catalogers have devoted much thought to who (or what) can be considered the author of a resource. Early cataloging codes kept things simple and emphasized recording the information that was present on the book, assuming that library users could work out for themselves whether a name included in the record was truly an author or some other role. Nineteenth-century cataloging theorist Charles Cutter defined the author simply as the writer of a book - the person who was responsible for its existence.

Later cataloging rules began to stretch the boundaries and add some nuance. A 1941 revision of ALA's cataloging rules provided guidance for assigning authorship to mediumistic writings (they recommended designating the medium as the author rather than the spirit). Another revision in 1949 extended the definition of author further to include roles such as artist, composer, photographer, and cartographer; while under the Anglo-American Cataloging Rules released in 1967, an author was "a person or corporate body chiefly responsible for the creation of the intellectual or artistic content of the work." The second edition of AACR, in 1978, reversed the earlier guidance around "spirit communications," and instructed catalogers to treat the *spirit* as the author, and the medium or other intermediary as a contributor. AACR2 also introduced standardized rules for formatting names in the LC Name Authority File, which is the database of established names for persons, corporations, families, and even titles of works. AACR2 included instructions for adding spirits to the LCNAF, but did not mention any other special categories of authors or beings.

With the implementation of Original RDA, several decades later in 2013, the definition

of author returned to being text-centric, but RDA also allowed for numerous other types of creatorship, such as compiler or remix artist. Additionally, RDA subtly expanded the definition of "person" to include "Christian saints, spirits, persons named in a sacred scripture or apocryphal book, fictitious or legendary person, or a real non-human entity" - since persons can be authors, by extension a being from any of these categories can be considered to have authored a text, can be contributed to the LCNAF, and can be added to a bibliographic record as an "author". We lived under this halcyon guidance for a little under a decade, and then the followup version of RDA, now known as Official RDA, snapped the definition of authorship back to "real human persons or collective agents [i.e. corporate bodies]". Under Official RDA, the relationship of fictitious characters to a resource can be recorded within a bibliographic record as a free text note (not a structured relationship). Other non-human entities can be recorded, imprecisely, as a "related entity" of a resource. The status of these categories within the LCNAF moving forward is in flux.

Can an AI be an author?

Cutter/ALA rules	sure, why not
AACR1	probably not
AACR2	probably not, but maybe?
Original RDA	yes!
Official RDA	absolutely not

Perhaps unsurprisingly, none of the existing cataloging standards directly address AI in any form (nor do more specialized standards such as OLAC Best Practices, MLA Best Practices, or the DCRM suite). Additionally, AI does not easily fall into of the categories of non-human entities I've mentioned so far. Given their treatments of other potential non-human authors, though, we can extrapolate how an AI attributed as the author of a work might be recorded following these various rules.

In present-day cataloging, AACR2, Original and Official RDA may all be used, depending on your institutional preference and any specialized guidelines you follow. So how does a modern-day cataloger navigate all these changing standards?

Cataloging of Resources Generated Using Artificial Intelligence (AI) Software **(PCC, February 2024)**

- treat AIs as software - i.e., not an author
- establish in the name authority file as a title
- include in catalog record as "related work"

example:

245 02 A book about artificial intelligence / Ꞥc written by ChatGPT
730 0_ Ꞥi Related work: Ꞥa ChatGPT

PCC to the rescue! In February of this year, the Program for Cooperative Cataloging released guidelines on "[Cataloging of Resources Generated Using Artificial Intelligence \(AI\) Software](#)," which makes recommendations for how to establish authorized forms of the names of AI software in the authority file, and how they should be entered bibliographic catalog records. PCC's guidance instructs catalogers to treat AI as software, and establish it in the name authority file as a title heading - that means that it is referred to by that title, and treated the way other types of software would be in the authority file, i.e. as a created work in its own right.

From the authority control perspective, this is fine and even good! It emphasizes that many things we refer to as "AI" these days are essentially just text generators operating via instructions laid out in code. ChatGPT, at its core, is just some lines of Python. (Okay, many lines of Python.)

However! This means that the main way an AI can be related to a created work is as a "related work". Which is not wrong, but it's generic, and in an era of cataloging where we can make note of very specific relationships such as "free translation of work" and "music for radio program work," it's a little unsatisfying to me. I'm also not sure it reflects how users would search for or think about a resource. I wanted to explore some other options, either as alternatives or additions to PCC practice, of representing an AI's relationship to a generated text.

are there other options?

- ~~don't record the information~~
- ~~ignore the PCC~~

So are there other options?

Don't record the information at all: since this talk is premised around the notion that we *do* want to record an acknowledged AI author, that's not an option.

Don't follow PCC rules: also not an option! If a cataloger is establishing names in the LC Name Authority File, they must follow PCC rules, full stop - I couldn't just decide add an AI to the LCNAF as a person or corporate body. Additionally, UMD is a PCC member institution, so we generally try to follow its guidelines for our bibliographic cataloging in addition to our authority control work.

re-use existing options

- use other authority files
- use existing MARC fields

567 __ †a ChatGPT generated the text †b ChatGPT †2 naf

567 __ †a DALL-E generated the pictures †b DALL-E 3 †2 wikidata

- explore other encoding standards

We could get creative with existing options.

Use another authority file: we prefer to use the LCNAF whenever possible, because it is robust, well-maintained, widely used, and is implemented by most ILS vendors (including our own new Alma). However, there are supplementary options. We could record an author from another authority files that has different standards, we could implement non-library options like Wikidata, or we could maintain a local authority file with special exceptions for AI programs. This would probably involve a fair amount of planning and configuration for relatively little gain.

We could explore other ways to record this information in the MARC format. This might be my favorite option - there are many lesser-used niche MARC fields out there, and it was fun to explore them. Fields in the 3xx and 5xx ranges are often used to record technical specifications, which seems promising: the [538 field](#), for instance, is used to record a note about "system details." However, the 538 and most other 5xx notes are all free-text rather than structured data. The 3xx fields are more structured, but none of them as currently established are suitable to provide information about software used to generate a work. The [567 field](#), though, offers a potential middle ground. It's used to record a "methodology note" consisting of "Information concerning significant methodological characteristics of the material, such as the algorithm, universe description, sampling procedures, classification, or validation characteristics." It contains a subfield †a for free text but also a subfield †b for controlled vocabularies, which could allow for both structure and nuance. The subfield †b supports use of both LCSH and LCNAF, as well as a number of other thesauri, which enables both flexibility and precision. This gets at the aspect of an AI's role as a

process, not just another work.

I've talked mainly about MARC, but there are other encoding standards in active development. BIBFRAME, for instance, has a property called [ProductionMethod](#) - it's currently used for physical production processes only, but could be modified to include digital production processes too. There are also a number of affinity groups developing "extension" ontologies for BIBFRAME - one could imagine an AI-focused group! Of course, you would still be limited by the descriptive standard you use.

propose new options

- propose (or just use) new relationship designators/labels
 - 730 o_ †i Generated by: †a ChatGPT
- propose a new MARC field (or BIBFRAME property, or RDA entity) entirely
 - MARC 731 Established Heading Linking Entry - Process Type
- create a new controlled vocabulary for AI processes
 - 657 _7 Gemini †2 aipt

We could propose new options that expand descriptive or encoding standards.

We could propose - or implement independently - new types of relationship designators or relationship labels. Original RDA allows for cataloger-created relationship terms if no existing term is adequate; however, Official RDA has a much more restricted set of relationship elements. While PCC has created a set of alternative relationship labels to the ones included in Official RDA, the list does not currently include any terms that adequately capture the idea of a work that itself generates another work, but we could potentially propose a new one. This would involve lots of planning and configuration for relatively little gain.

We could also propose an entirely new MARC field or subfield to support the recording of information about software used to create a resource, which could also be useful for other types of resources such as cartographic materials and datasets (and would be popular with scholars studying the history of word processing software, although that would require a lot of retrospective work).

We could also create an entirely new vocabulary specifically for AI processes: similar to the existing [LC Medium of Performance Thesaurus](#) or [LC Demographic Terms Thesaurus](#), we could have an "AI Process Thesaurus." Its vocabulary terms could be recorded in the existing [MARC 657 field](#), "Index Term - Function," which records a "Term specifying the function or activity which generated the materials" - the 657 currently has only one thesaurus allowed for use, but we could propose this new thesaurus be added.

Thank you! Questions?

So, for my starting questions - *Who wrote this?* and *Who is an author?* - I don't have any fully satisfying answers, and I suspect that these issues will remain an area of lively debate for some time to come. However, I hope I've been able to propose some "innovative practices" for the representation of AI resources in the catalog!