1  **Name Authority Control in Repositories**

2  **By**

3  **Charity K.M. Stokes, Cataloging and Metadata Librarian, Texas A&M University Libraries**

4  **David B. Lowe, Digital Collections Management Librarian, Texas A&M University Libraries**

5

6  *Introduction*

7  When users search for information on a particular topic or for works by a particular person,
8  there are certain assumptions they make. One of the most cherished is that they can easily find
9  all works on or by a given person with minimal effort and, if they have the name wrong, that
10  they will be pointed in the right direction as to what form of that person's name to use in their
11  search.  However, neither this assumption nor the functionality of collections comes about
12  naturally.  It is a result of how the metadata is entered and organized.  One of the guiding
13  principles of information organization is that of authority control.

14  What is authority control?  It is a way to gather all variations of the names of a person,
15  corporate entity, or subject into one authorized access point. Someone could be known by a
16  nickname, their initials, a stage name, or a pseudonym, yet each variation refers to the same
17  person.

18  Most institutional repositories (IRs) are found in academic institutions where the primary
19  resource is the tenured faculty members whose critical career need is to be able to chronicle
20  their work.  As the push for open access gains steam, one of the ways faculty can show the

21  impact of their work is simply by how many downloads, references, and views their work has

22  had, which IRs can provide.  But a key prerequisite is that faculty need to have their work
23  represented under one name.  If their name appears in multiple ways, it is a challenge to gather
24  all of that information under one heading.

25  *Pre-Repository Authority Control*

26  One of the most interesting aspects of modern authority control is that it did not originate from
27  user needs.  It began simply as a tool, an in-house file used by catalogers so that they could
28  bring all the works of a single person under an "authorized" name.  The catalogers created the
29  preferred heading – usually garnered from the first book by that person added to the local
30  collection.  Any other works by that author were then placed under the established heading for
31  that library's collection.

32  So, in a quirk of library history, what is now a major access point started out in local technical
33  service departments as a part of the catalogers' toolbox.  Charles Cutter, the renowned creator
34  of cataloging rules in the U.S., even stated that authority control was to be performed for the
35  convenience of the cataloger (Cutter, 1891).  In the grand old days of card catalogs, the

36 authority file was kept in the back room of the cataloging unit.  A new book would come in, the
37 cataloger would check the author listing within the book against the local authority file and
38 then create the bibliographic card under the heading of the author utilizing the locally
39 authorized name.

40 When Cutter wrote his rules, collections were a great deal smaller, perhaps consisting of a
41 thousand or so books. In the early 20th century, most libraries' collections included around
42 3,000 volumes per capita (Kevane & Sundstrom, 2014). In 2012, the average size of collection in
43 a United States public library was 110,708 items with a median of 42,833 (Grimes, Manjarrez, Miller,
44 Owens, & Swan, 2014).  Because of the smaller size of collections, authority control was useful to
45 the cataloger, but not obvious to the user given that few authors would share the same name
46 as another.  However, as book publication increased, so did the possibility that two or more
47 authors would share the same name, increasing the need for authority control.  However, such
48 work is labor intensive, which limited the extent of its implementation.

49 In the 1930's, to address the growing need, many libraries began depending on the Library of
50 Congress (LC) for cataloging cards.  As part of this service, LC also began making cards available
51 for the local authority files.  Local name authority files began to reflect the same authorized
52 name headings found in other libraries across the country.  It did not take long before the
53 author file cards became the accepted form of an author's name.  Thus, a national name
54 authority file was born.

55 Still, cataloging manuals at the time paid scant attention to *how* to do authority work.  Cutter's
56 Rules only say: "Give the names, both family and Christian, in the vernacular form, if any
57 instance occurs of the use of that form in the printed publications of the author" and "when an
58 author's name is various spelled, select the best authorized form as heading, add the variants in
59 parentheses, and make references from them to the form adopted"  (Cutter, 1891) (pp 24-25).
60 Even the original Anglo-American Cataloging Rules (AACR), published in the 1960's, only advised
61 to "make a heading under author's name in full and in vernacular form … enter under family
62 name followed by forenames and dates of birth and death for specific identification when
63 available" (Wynar, Tannenbaum, & Christensen, 1966).  Neither explicitly states how to create a
64 separate authority file, only that the cataloger is to refer to the authorized name.  Catalogers
65 developed authority files as their primary tool to know what the authorized form of names
66 (personal and corporate) and subjects would be (Auld, 1982).

67 It was only with AACR2 (1978) that approved and standardized authority control practices came
68 to be (which included an entire chapter on the *see* and *see also* references) (American Library
69 Association, 1967). However, once again, the emphasis was on the fact that references *should*
70 be made, but now *how* to do it. In 1978, *Authorities, a MARC Format* was published. It set a
71 national baseline for automated authority records, based on the American National Standards
72 Institute (ANSI) standard for the communication of authority records by means of magnetic
73 tapes.  This preliminary guidance only carried an implicit standard for quality, "with
74 specifications and content designators for name, uniform title, and subject authorities,"  (Auld,

75    1982) (p. 323) thus establishing the differing types of electronic authority files we have today.

76    This preliminary edition was replaced in 1981 when LC published the first edition of *Authorities,*

77    which also included the addition of series authority and series treatment. With this publication,

78    a national standard became available for the recording, structuring and sharing of authorities

79    for names (personal and corporate, uniform titles, subjects and series).

80    Another development that was driving the push for clear standards and rules was the advent of

81    the Online Public Access Catalog (OPAC).  When Machine Readable Cataloging (MARC)

82    appeared as a bibliographic standard in 1968, the same need for standardized names that was

83    first felt in the infancy of modern cataloging now required even more attention as the

84    quantities scaled upward.  Users still needed to identify authors and the complete listings of

85    their works.  The only thing that changed was the scope of the bibliographic universe.  No

86    longer was it just about what the local collection contained, but the entire library community.

87    It was not just libraries who were concerned about authority control; the publishing industry

88    was no less affected. Bowker was one of the very first to compile an authorized list of authors.

89    With their iconic publication *Books in Print*, Bowker also needed to know that the books written

90    by one author were attributed to that author.  So, in 1981 Bowker published *Authors' Names;*

91    *An Authoritative Listing of Personal and Corporate Names,* which was based on LC records

92    (Bowker, 1981).

93    It was the development of the computer age that accelerated these developments. Catalogers

94    began sharing their knowledge as well as their records via cataloging utilities. The Library of

95    Congress Authority File (LCNAF) became the definitive 'authority file' for the country and most

96    of the western hemisphere. By using a standardized and trusted source, libraries reduced the

97    overhead for cataloging by doing away with local authority files.  It increased productivity and

98    reduced the cost associated with cataloging.

99    Authority control changed tremendously with the introduction of *Functional Requirements for*

100   *Bibliographic Records (FRBR)* (International Federation of Library Associations, Study Group on

101   the Functional Requirements for Bibliographic Records., 1998)  and the *Functional*

102   *Requirements for Authority Data (FRAD)* (IFLA Working group on Functional Requirements and

103   Number of Authority Records (FRANAR), 2008)*. This was a radical shift away from Cutter's

104   requirements to identify and disambiguate objects of a catalog to fulfilling the specific FRBR

105   user's tasks.  The tasks are

106   •   To **find** entities that correspond to the user's stated search criteria

107   •   To **identify** an entity

108   •   To  **select** an entity that is appropriate to the user's needs

109   •   To acquire or **obtain** access to the entity described

110

111   ***Within Repositories***

112    To relate this historical narrative with IRs, consider that, with digital collections, the operative
113    word is 'collections'.  A library's collections need to be accessible.  Just as libraries provided
114    access through book indexes, card catalogs, OPACs, etc., metadata librarians create a surrogate
115    record (metadata) in order for users to locate the information they seek.  One of the most
116    important aspects of what users need – especially academic users – is to find all the works by
117    an author, or any works that a person has contributed to.  To this end, the metadata for
118    authors (and subjects and series) needs to be <u>collocated</u> – gathered under a consistent,
119    authorized form. The tools we use may have changed, but the needs of our users have not.

120    Focusing on the IR context, a worthwhile preliminary discussion might include the evolution of
121    the terms "digital collections" and "institutional repository" and the library functions associated
122    with them—such as they have come to be—since the advent of the Web.  Initially, digital
123    collections as a phenomenon in academic libraries grew largely out of in-house scanning
124    operations that pre-date networked information tools like web browsers.  The content sources
125    were holdings from the libraries' archival and circulating collections:  photographs, postcards,
126    maps, and most certainly print items such as books and journals, but also audiovisual materials
127    like local oral histories.  Such digital conversion activities had begun in libraries before networks
128    blossomed in the 1990s, with file sharing happening via the various evolving media of the time.
129    After the release of the first graphical tool Mosaic in 1993, browsers became the obvious frame
130    for all of this content.  At the same time, capture equipment to produce that digital content
131    from analog objects (like books and paper) became more affordable, which led to more content
132    getting created and available.  The following decade saw the rise of mass digitization efforts,
133    particularly of textual objects, such as those led by the Internet Archive, Microsoft, and the
134    Google Books project, parts of all of which have been collocated spectacularly in HathiTrust.

135    Concurrently with this conversion from paper, responsible information professionals also
136    devoted their attention to born digital materials.  Purists would argue that a true IR consists
137    primarily of "born digital" materials, predominately electronic theses and dissertations (ETDs),
138    articles, conference papers, and presentations; "digital collections," on the other hand, they use
139    to refer more to sets of items digitized from analog originals, perhaps mixed with some born
140    digital materials.  The first university to require electronic submission of theses and
141    dissertations was Virginia Tech, in 1997 (https://vtechworks.lib.vt.edu/handle/10919/5534 ),
142    but the true arrival of IRs may be comfortably dated to the initial development and release of
143    the DSpace software in 2002-2003, mentioned by Clifford Lynch in his seminal work advocating
144    for IRs. (Lynch, 2003)

145    One feature of IRs that complicates the metadata aspect is that much of the ingested content
146    (OA articles, gray literature) was not as routine to cataloging workflows and the staff that ran
147    them.  If they were even involved at all with the new IR materials workflows, typical cataloging
148    departments at the latest turn of the century were heavily involved in monograph, serial, and
149    A/V materials cataloging, but were less accustomed to material without a book-like title page or
150    "chief source of information,"—as the trade terminology goes--which typically would have

151 fallen to more specialized and experienced "original catalogers," as they are known.  The sheer
152 numbers of these new documents, the lack of formal identification of their authors, and
153 unmediated deposit (with authors or their designates left to their own devices), has led to an
154 accumulation of content in IRs that seems utterly devoid of traditional authority control.

155 It was considered as an advantage that the authors could input their own metadata.
156 Structurally, the idea was still the same:  describe the work and provide access points.  And who
157 could better describe a work than the authors themselves?  Having authors or their designates
158 essentially catalog their own works would cut out the middle man and allow faster access to
159 material.

160 Such a tack was functional for many years.  Basic access fields such as author and title were
161 present. However, because there was no authority file, variations in names began to crop up. At
162 first, the variations were a minor inconvenience.  But just as with physical collections, as digital
163 collections grew, so did the tangles created by the lack of authority control.  Users began to
164 grow frustrated trying to figure which of the five Professors J. Vance were they looking for.
165 Was it James Vance?  Or Joan Vance?  Jack, Jill, or John?  And if it was a James Vance, did that
166 include James A. Vance and was he the same as James Allen Vance? Which one was the
167 composer and which was a respected professor of biology? Users became frustrated and would
168 often would give up when trying to find the works they were looking for by a given author.
169 There were variations due to the fact that occasionally authors would use their full names,
170 while at other times their initials.  And then there is the issue simply of having an extra space or
171 a misspelling, creating endless variants.

172 ***Challenges***
173
174 A suitable review of the topical literature would be a summary of the challenges reported in the
175 scholarly record related to authority control in IRs.  The relatively brief history of IRs, in tandem
176 with the concentration of early related work being heavily devoted to advocacy for IRs as a
177 concept, translates into a quick exercise.  Among the first to call attention to the issue of
178 authority control in IRs in a comprehensive, studied manner was Salo (Salo, 2009) who outlines
179 problems with the available tools and related workflows but also speculates about some fruitful
180 paths toward resolutions.
181
182 One issue brought up by Billey  (Billey, 2019) has to do with privacy.  Current standards for the
183 industry standard Name Authority Cooperative Program (NACO) authority records are based on
184 FRAD, which greatly expanded the number of attributes to describe people.  Under Cutter's
185 rules, the only need for attributes was to identify and disambiguate the names in order to
186 facilitate discovery. However, when FRAD was codified into Resource Description and Access
187 (RDA) in 2010 and established as the standard through the RDA Toolkit (Joint Steering
188 Committee for Development of RDA., Chartered Institute of Library and Information
189 Professionals., American Library Association., & Canadian Library Association., 2010), catalogers

190 went from disambiguation to actually describing people.  This descriptive process included not
191 just information that was already in use (name of person, dates such as year of birth and/or
192 death, fuller form of name) but to much more detailed information such as:  gender, country,
193 address, profession, titles, and affiliations.
194
195 On the surface, many US based academic institutions would see no problem with including the
196 above information.  However, given the current political atmosphere, something seemingly
197 innocuous as affiliation can cause personal problems.  For example, perhaps there is a popular
198 children's author who is member of a minority religion.  This personal affiliation has nothing to
199 do with their writings or works, but can result in them being targeted by the more radical
200 elements that exist in society.  Then there is the issue of gender.  If the person's gender –
201 especially if non-binary – impacts their work, it will be obvious in the work itself.  If not, it really
202 does not serve any purpose in the identification or disambiguation of the author's name.  So
203 the challenge for authority control in repositories is the same as for many online media – that
204 of balancing the privacy of the individual against the needs of the organization.
205
206 A final challenge faced by repositories is pragmatic.  Budgetary constraints affect the ability to
207 manage metadata by limiting the number of resources available to perform that management.
208 It has long been known that authority control by people is an expensive investment in
209 bibliographic control. Shrinking technical service departments and the outsourcing of those
210 tasks has led to a dearth of expertise and a shrinking number of people able to perform the
211 work.  However, one possible solution is a type of human-machine hybrid system, leveraging
212 software to help control costs (Liu & Qin, 2014).
213
214 ***Current State of Name Control***
215
216 The authors' knowledge of practices springs from their experiences at a variety of institutions,
217 combined with a focused search for IR workflow instructions across many institutions in 2019. It
218 confirms that metadata for ETDs are typically handled manually by staff.   When it comes to the
219 names of authors, advisors, and committee members, the staff normalize names by referencing
220 standard authorized forms (*e.g*., LCNAF).  If a name does not yet have an authority record, the
221 staff usually follow a standard algorithm, such as "LastName, FirstName MiddleInitial."  Sadly,
222 even the best of algorithms can lead to conflicts, due to ambiguities introduced with features

223 such as compound last names or life events (such as marriage) that result in name changes.

224 Less mediated workflows, such as self-deposited articles, manuscripts, or presentations are
225 even more prone to conflicting entries.  External tools such as OpenRefine (Carlson & Seely,
226 2017) may be enlisted to help resolve existing entries, while some IR software has limited
227 internal functionality in this arena. For example, EPrints (EPrints Project, n.d.); (Salo, 2009) and
228 DigitalCommons (Edwards, 2018) include the ability to merge name records.
229
230 Some authority control solutions are arising at the institutional and regional consortium levels.

231   Digital Library efforts at the University of North Texas (UNT) have produced the UNT Names
232   App (https://digital2.library.unt.edu/name/ ).  Covering primarily personal and organizational
233   names, the application is incorporated into their IR workflows, which feature mediated deposit.
234   Similarly, librarians at the University of Houston (UH) Libraries have implemented an instance of
235   iQvoc called Cedar that covers the names of individuals and organizations, as well as subject
236   terms.  At Columbia University Libraries, one notable function of their metadata editing tool
237   Hyacinth is its capacity to mint URIs for named entities that lack them elsewhere.  A consortial
238   project at the Mountain West Digital Library (MWDL), known as the Western Name Authority
239   File (WNAF), seeks to create a central file for its partners that will be compatible with Linked
240   Open Data (LOD) efforts.
241
242   ***Future Trends: LOD***

243   Although LOD has not been mentioned heretofore in this chapter, it did not just pop up by
244   happenstance.  While having a name authority control system in place is a worthy, practical
245   cause in and of itself, in fact there are much larger implications.  Name authority systems that
246   produce and manage uniform resource identifiers (URIs) can relate their efforts to LOD
247   developments; in turn, the LOD efforts are connected to Artificial Intelligence (AI) and Machine
248   Learning (ML)—some  of the most promising and awe-inspiring accomplishments of the current
249   age, as components of autonomous vehicles and smart speakers, just to name some examples.
250   At the heart of structuring these advances for further success are open standards such as
251   Resource Description Framework (RDF) which relies on a basic tripartite grammar structure
252   "Subject/Predicate/Object" for which the subjects (or actors) and the objects (or those acted
253   upon) are most commonly represented by URIs.  No less than Sir Tim Berners-Lee (2006), in his
254   seminal and foundational piece on Linked Data, states expectation #1 as:  "Use URIs as names
255   for things."  So the austere world of authority control from librarianship brings us face to face
256   with the future of humankind's interaction with information technology in our everyday lives.

257   To be coldly realistic, though, such a future is fraught with huge risks where privacy and
258   personal security are concerned.  Insofar as librarians can collectively influence the directions at
259   hand, it is worth considering what is at stake.  As mentioned previously, Billey (2019) shares
260   sober advice about the vulnerable zone where authority work could impinge upon Personally
261   Identifiable Information (PII).  Referring to particular date, gender, and affiliation data points,
262   Billey cautions: "Recording this information could violate a person's privacy, make their
263   personal information vulnerable to bad actors, and even possibly put someone in danger" (pp.
264   10-11).  The solution offered is to be circumspect in what data points even get scoped, much
265   less recorded, in authority control systems.

266   Having considered the promises and risks for authority control systems related to local IR
267   implementations, it bears emphasizing that local systems will tend to be the most relevant and
268   familiar with the institutional context.  Open URIs that unfold thorough, accurate, yet not overly
269   revealing pockets of information about people and their groupings will have the capacity to

270 integrate IRs with the greater world of scholarly communication and facilitate positive
271 interaction.  Standardization will be critical, but at the moment those standards are not yet
272 established above any threshold of collective refinement.  The near term ahead will be full of
273 experimentation and surely some solid best practices will emerge.

**Conclusion**

275 For the user, the need to be able to identify all the works by a given author and know that this
276 author is the correct one has not changed since Cutter's day.  What has changed is the sheer
277 number of authors available in the bibliographic universe.  For institutional repositories, this
278 includes not just the traditional book author, but also authors of articles and grey literature.
279 The scale is overwhelming for the metadata specialist.  For a user, it can easily go beyond
280 overwhelming to baffling. This chapter covered the development of authority control as well as
281 the constraints and challenges inherent in trying to impose authority control.

282 Authority control is a needed tool for our repositories and digital collections.  Given the
283 emphasis in academia on citations and analytics for the purposes of career advancement,
284 faculty need to cite metrics connected to how often their works were taken up in their
285 professional communities in order to convey impact.  Graduate students would like to know
286 what disciplines a professor has published in, as well as which professor served as advisor for
287 other graduate students' works.  No one wants to wade through multiple publications for one
288 person who appears with as many as 12 variant names!  So much needed information is lost in
289 such a welter of name variants.

290 The rapidly developing semantic web enables a world where users can gather the information
291 they need by following the connections between different entities – whether a person, a group,
292 or a subject.  The vision of linked data is the foundation upon which current standards (RDA,
293 FRBR, FRAD) have been built to work in a world of artificial intelligence and machine learning.

294 Our world is in the middle of a revolution, an information revolution that is no less a seismic
295 shift than the industrial or technological revolutions before it. Digital collections and
296 repositories are not only a product of that revolution but are helping to drive it.  Authority
297 control will be part of the steering.

298 ## References

299 American Library Association. (1967). *Anglo-American cataloging rules: North American Text.* (C. S.
300     Spalding, Ed.) Chicago, IL: American Library Association.

301 Auld, L. (1982). Authority Control: an eighty-year review. *Library Resources & Technical Services*, 319-
302     330.

303 Berners-Lee, Tim.  (2006). Linked Data.  *Design Issues.*  Retrieved from:
304     https://www.w3.org/DesignIssues/LinkedData.html

305    Billey, A. (2019). Just because we can, doesn't mean we should: An argument for simplicity and data
306         privacy with name authority work in the linked data environment. *Journal of Library Metadata*.
307         doi:10.1080/19386389.2019.1589684

308    Bowker. (1981). *Authors' names: An authoritative listing of personal and corporate names.* New York:
309         Bowker.

310    Carlson, S., & Seely, A. (2017). Using OpenRefine's reconciliation to validate local authority headings.
311         *Cataloging and Classification Quarterly, 55*(1), 1-11. Retrieved from
312         https://doi.org/10.1080/01639374.2016.1245693

313    Cutter, C. A. (1891). *Rules for a dictionary catalogue.* Washington: Government Printing Office.

314    Edwards, L. (2018, June). Authority control in digital commons: Why bother? *Presented at the Digital*
315         *Commons Southeatern User Group Meeting*. Johnson City, Tennessee. Retrieved from
316         https://encompass.eku.edu/fs_research/280

317    EPrints Project (n.d.).  *Authority lists - EPrints Documentation*.  Retrieved from
318          https://wiki.eprints.org/w/Authority_Lists

319    Grimes, J., Manjarrez, C., Miller, K., Owens, T., & Swan, D. W. (2014). *Public Libraries in the United States*
320         *survey: Fiscal Year 2012.* Washington: Institute of Museum and Library Services.

321    IFLA Working group on Functional Requirements and Number of Authority Records (FRANAR). (2008).
322         *Funcational Requirements for Authority Data: A conceptual model.* München: KG Saur. Retrieved
323         from https://www.ifla.org/files/assets/cataloguing/frad/frad_2013.pdf

324    International Federation of Library Associations, Study Group on the Functional Requirements for
325         Bibliographic Records. (1998). *Funcational requirements for bibliographic records: Final report.*
326         Munchen: K.G. Saur. Retrieved from https://www.ifla.org/files/assets/cataloguing/frbr/frbr.pdf

327    Joint Steering Committee for Development of RDA., Chartered Institute of Library and Information
328         Professionals., American Library Association., & Canadian Library Association. (2010). *RDA*
329         *Toolkit.* [Chicago}, IL: American Library Association.

330    Kevane, M., & Sundstrom, W. A. (2014). The development of public libraries in the United States, 1870-
331         7930: A quantitative assessment. *Information & Culture, 49*(2), 1117-144.

332    Liu, Z., & Qin, J. (2014). An interactive metadata model for structural, descriptive, and referential
333         representation of scholarly output. *Journal of the Association for Information Science and*
334         *Technology*, 964-983.

335    Lynch, Clifford A.  (2003). Institutional Repositories:  Essential Infrastructure for Scholarship in the Digital
336         Age.  *ARL:  A bimonthly report on research library issues and actions from ARL, CNI, and SPARC.*
337         Retrieved from:  https://www.cni.org/wp-content/uploads/2003/02/arl-br-226-Lynch-IRs-
338         2003.pdf.

339    Salo, D. (2009). Name authority control in institutional repositories. *Cataloging & Classification*
340         *Quarterly*, 249-261. doi:10.1080/0613937092737232

341    Wynar, B. S., Tannenbaum, E., & Christensen, C. (1966). *Introduction to cataloging and classification: A*
342        *teaching guide with illustrations of major principles for descrptive cataloging and classification*
343        (Second, revised and enlarged ed.). Denver, CO: Colorado Bibliographic Institute.

344

345