



Being Irrelevant: Why Library Data Interchange Standards Have Kept Us Off the Internet

Kenning Arlitsch

This is a preprint of an article that originally appeared in the Journal of Library Administration in 2014.

Kenning Arlitsch, (2014) "Being Irrelevant: Why Library Data Interchange Standards Have Kept Us Off the Internet", Library Hi Tech, Vol. 54 Iss: 7, pp.609-619 <http://www.emeraldinsight.com/doi/full/10.1080/01930826.2014.964031>

Made available through Montana State University's [ScholarWorks](http://scholarworks.montana.edu)
scholarworks.montana.edu

Column Title: posIT

Column Editor: Kenning Arlitsch, Dean of the Library, Montana State University, Bozeman, MT
kenning.arlitsch@montana.edu

This JLA column posits that academic libraries and their services are dominated by information technologies, and that the success of librarians and professional staff is contingent on their ability to thrive in this technology-rich environment. The column will appear in odd-numbered issues of the journal, and will delve into all aspects of library-related information technologies and knowledge management used to connect users to information resources, including data preparation, discovery, delivery and preservation. Prospective authors are invited to submit articles for this column to the editor at kenning.arlitsch@montana.edu

Being Irrelevant: Why Library Data Interchange Standards Have Kept Us Off the Internet

By Kenning Arlitsch

Introduction

I used Uber for the first time recently, and the experience was so superior that it will not be the last. Uber Technologies, Inc. has been described by its founder as a “technology platform that connects riders and drivers,” and the result is a four-year old company whose value now exceeds \$18 billion (Segall, 2014) and has millions of satisfied customers. Uber has been so successful for a simple reason: it conveniently delivers a product that people need. It uses existing and ubiquitous information delivery platforms (smartphones) that make it convenient to use, and it takes advantage of common data interchange technologies like GPS that allow a rider to know how soon and where the driver will appear. While Uber has its detractors – largely because it threatens the taxi industry - its meteoric rise offers an important technological lesson for libraries.

Conversations about the future of libraries invariably raise questions about “relevance.” One way to define relevance is to evaluate how well library “products” integrate into the popular information ecosystem, i.e. the Internet. It is in this ecosystem that Uber has excelled and libraries have struggled. To use library products our customers must deliberately move into another information ecosystem built by libraries and library vendors, when they should be able to discover and have seamless access in the ecosystem where they already conduct their business. Libraries force customers to use technological tools to which they are not accustomed, which in turn spawns an instruction mini industry.

Lorcan Dempsey coined the phrase “inside-out library,” indicating that discovery in the networked world happens outside the localized ecosystem of the library, and that libraries are challenged to push their products to the outside world (Dempsey, 2010). He and co-authors further explain that “effective discovery means syndication to search engines, to disciplinary resources, or to other specialist network-level resources” (Dempsey, Malpas, & Lavoie, 2014). In other words, libraries must deliver their products to the ecosystem where users live, work, and play, taking advantage of popular standards and platforms. This article speaks to the futility of trying to fulfill that imperative through data interchange protocols or publishing platforms whose use is almost non-existent outside the library world. Internet search engines, the true discovery systems of the Internet, have little or no use for the metadata standards and the data interchange protocols that libraries and archives have developed, and I question the wisdom of continuing to sink resources into using them to describe our resources

Background

“The protocols and standards that underpin our digital lives are, in the long term, as ephemeral as the sand mandelas of the monks” (Weibel, 2010).

The greater information ecosystem where people live and work is the Internet, and in that world the level of interoperability and integration has rapidly and steadily improved. Even major proprietary platforms don’t get in the way of data interchange much, anymore. Nearly every other business that trades in information or sells products has figured out how to leverage the protocols and platforms supported on the Internet. The number and variety of information, productivity or entertainment applications that people can access through the Internet search engines is stunning, but the vast collections owned by libraries barely register.

I don’t mean to say that library products are completely inaccessible. But the richness and accuracy of the descriptions in which librarians and archivists have invested countless hours don’t transfer well to the Internet, making our products far less discoverable and usable than they could be. We have had a tendency to implement standards and protocols that are indigenous to the library environment and are used nowhere else. Worse, those standards and protocols are often implemented through laborious processes that don’t scale to the data deluged world we now live in. MARC, EAD, Dublin Core, TEI, and OAI-PMH, are all examples of library metadata standards and data sharing protocols that have little or no impact outside libraries. Some progress is being made toward evolving these standards for the Semantic Web, but it is slow and I expect our workflow status quo to continue for years.

Internet search engines bias my perspective. After a decade of building digital library repositories I came to the conclusion that it was hard to justify the investment in digitization because use was so low. Research on search engine optimization (SEO) followed, and gradually my colleagues and I figured out how to increase visitation and downloads by improving the rate at which our digital objects

are harvested and indexed by search engines (Arlitsch & O'Brien, 2013). More recently, our research team has become interested in the role of the Semantic Web in SEO and the potential it offers for more accurate and descriptive discovery and rich use of digital repositories. The SEO research taught me that while there are multiple facets to this issue, the problems libraries face are mainly related to data structure and accuracy, as well as the ability to deliver seamless (aka "convenient") access.

"What we've got here is failure to communicate"

The line from the 1967 film *Cool Hand Luke* starring Paul Newman is apropos for this discussion. The fact that I can type that quote into Google and immediately find a link to a YouTube clip that plays the exact scene in the film reveals the comparative poverty of what digital libraries can offer users in those same search engines. What are the chances that hits from deep inside archival collections will appear in an Internet search, and if they do appear, will they lead so quickly to a digitized photograph or video that displays right in the browser? Search engines want to deliver the best experience possible to their users, and burdensome or inconvenient experiences are unlikely to result in referrals.

Standards inconsistency

"Standards are great – everyone should have one" (Group, Systems, Institute, Forum, & International, 1991)

The humor in Robert Metcalfe's famous quote belies a serious problem with library efforts to share content with the world. It's not only that libraries have developed data interchange standards that search engines have no interest in, but it's also the inconsistency with which those standards are applied. Librarians sometimes have trouble agreeing on definitions of specific fields, which may help explain how Dublin Core became a metadata standard of the lowest common denominator. Archivists have similar difficulties with the definition of EAD fields, a problem exacerbated by the propensity of some archivists to claim uniqueness among their peer institutions: "we do it differently here."

Manual re-keying of metadata adds to problems of accuracy and consistency. Librarians, archivists and paraprofessional staff do a lot of manual keying as they prepare bibliographic or descriptive records for library catalogs, finding aids, digital collections, or institutional repositories. Studies of medical transcriptions have found error rates as high as 46% (Seeley, Nicewander, Page, & Dysert, 2004), and we should assume that transcribing introduces errors into library metadata, affecting machine readability and thus search engine ingest. Although controlled vocabularies have existed for decades, the slow processes surrounding approval of new terms, maintenance, and integration tools have led to a proliferation of unofficial terms. Expanding linked data sources offer a more automated method of importing terms and names, but most library software doesn't yet support that possibility.

Cataloging and Metadata

Following is a brief discussion of various data interchange efforts that have evolved in libraries and archives.

Machine Readable Cataloging (MARC)

First discussed in the late 1950's and formally launched in the mid-1960's, MARC may be considered the grandfather of library automation standards. The Library of Congress led development of the standard and by 1974 it reported 74 institutional subscribers to its MARC cataloging service (Avram, 1975). It took several more decades for MARC to gain full adoption by libraries, and for automated catalog systems that could support MARC to grow and mature. The MARC standard has enabled the creation of many millions of structured data records, but "this data is not readily available for re-use outside of the library community" (Styles, Ayers, & Shabir, 2008).

Even before search engines came into being the development of the Web in the mid 1990's quickly demonstrated the shortcomings of MARC for this new environment. "As the Web exploded on our desktops, it was evident that MARC cataloging of electronic resources would be too complex and costly, and might not be suitable for many electronic resources in any case" (Weibel, 2010). OCLC, probably the world's largest aggregator of MARC records has stated that "Bibliographic data stored in traditional record formats has reached its limits of efficiency and utility." (OCLC, Inc., n.d.), and is pushing its linked data research and development efforts on that basis. While powerful for its day, the MARC framework was based on text strings and although conversion to RDF is in the works it has not evolved quickly enough to an entity-based structure to which URIs can be assigned "to refer to the resource that we, as people, infer from the literal string" (Styles et al., 2008). In other words, people are smart and can infer meaning from text strings. Machines have trouble making those inferences and need authoritatively established entities to help them understand information.

Text Encoding Initiative (TEI)

Development of TEI began in 1987 to address "the proliferation of systems for representing textual material on computers" (Mylonas & Renear, 1999). The ambitious effort led to a robust data description language for the humanities whose first official version was released in 1994, and has seen considerable application in some research libraries and rare books archives. Like most encoding standards it can be applied to differing levels of complexity, requiring varying degrees of labor.

In their article celebrating the ten-year anniversary of the birth of TEI, Mylonas and Renear discuss a condition in academia that reveals a microcosm of the current problem of data interchange across strata: "Not only do different disciplines have quite different interests and perspectives, but also, it seems, different conceptual schemes: fundamentally different ways of dividing up the world. What is an object of critical contest and debate for one discipline is theory-neutral data for another, and then completely invisible to a third" (Mylonas & Renear, 1999). This condition

is currently magnified in scientific disciplines, as researchers struggle to make their data sets discoverable and usable across disciplines. Interdisciplinary research holds great promise for new scientific discoveries, but this promise is hampered by monolithic and discipline-specific ontologies. The National Institutes of Health recognize this as a critical shortcoming that hampers new developments in research and have accordingly begun funding “Big Data to Knowledge (BD2K)” initiatives (National Institutes of Health, n.d.). The notion that libraries might be able to help with this initiative seems far-fetched.

In the early days of TEI some enthusiastic supporters suggested that HTML was insufficient for the growing Web, and that the TEI offered a model for numerous document type definitions (DTD) that could be created and managed through SGML (Barnard, Burnard, DeRose, Durand, & Sperberg-McQueen, 1995). They were right about HTML, but TEI itself has remained insular and unable to deliver its own rich data to the Internet.

Encoded Archival Description (EAD)

The alpha version of EAD was released in February 1996 and has since been widely adopted in the archival community. Developed to bring standardization and machine readability to the painstaking detail that archivists create in finding aids, EAD promised to “support the long cherished dream of providing archivists and both professional and public researchers universal, union access to primary resources” (Pitti, 1999).

Search engines were in their infancy in the late 1990’s, and Google didn’t establish its dominance until the early 2000’s, but even as EAD evolved from its original instantiation in Standard Generalized Markup Language (SGML) to the less complex and more flexible Extensible Markup Language (XML) there was cautious optimism of discovery through the Web. “The use of XML may soon solve this problem, as all XML/EAD finding aids will be fully searchable across the Web, but, of course, we do not yet know how XML and new metadata RDF (Resource Description Framework) standards will affect retrieval, nor when Web search engines will accommodate XML” (Tibbo & Meho, Lokman I., 2001). That optimism has unfortunately not been rewarded, mainly because Internet search engines seem to have no interest in EAD. Efforts to make finding aids discoverable through the web have ranged from “dumbing down” a subset of EAD fields to Dublin Core metadata for discoverability alongside digital collections, to rendering them as HTML, to converting a subset of fields to MARC. OCLC’s ArchiveGrid is probably the best known example of this latter method and has collected an impressive three million records for finding aids (Washburn, Eckert, & Proffitt, 2014). While this effort reveals bibliographic data for finding aids the richness of EAD is lost in the conversion, just as it is with the prior methods that utilize Dublin Core or HTML.

The CENDARI project represents a step in the right direction, as it recognizes the inherent interoperability difficulties of EAD and stresses the need for semantic application and machine analysis: “...archival descriptions should act as datasets,

which can be subject to detailed analysis...capable of supporting machine based techniques used more commonly for analyzing large datasets or text corpora.” (Gartner, 2014) But CENDARI falls short by relying on community-developed ontologies that are admittedly lacking at the moment.

Dublin Core

The Dublin Core Metadata Initiative (DCMI) was born in the mid 1990’s to address the difficulty of finding resources on the Web. Originally a set of 15 elements, Dublin Core (DC) has been enhanced to include additional qualifiers and it has become the de facto metadata standard for digital collections and many institutional repositories. The evidence as to whether Google makes use of DC metadata that is pulled from digital asset management databases and embedded into HTML display pages is mixed (Cutts, 2009), but Google Scholar advises against its use for academic papers: “Use Dublin Core tags (e.g., DC.title) as a last resort - they work poorly for journal papers because Dublin Core doesn’t have unambiguous fields for journal title, volume, issue, and page numbers” (Google Scholar, 2011). Instead of Dublin Core, Google Scholar recommends using one of four non-library metadata schemas: Highwire Press, PRISM, Eprints, and BePress (Arlitsch & O’Brien, 2012).

Open Archives Initiative Protocol for Metadata Harvesting (OAI-PMH)

OAI-PMH was developed in 2000 and grew out of a “vision to stimulate the growth of open e-print repositories...” (Lagoze & Van de Sompel, 2003). It allows repositories to offer Dublin Core metadata in an XML bitstream to harvesters that honor the protocol, and it may be the most common method in libraries to harvest and aggregate metadata. Many library vendors also support the protocol for their digital asset management software and discovery systems. However, compliance is questionable. Richard Urban found that less than 25% of IMLS National Leadership grant projects were providing item-level metadata using OAI-PMH in 2011 (Urban, 2011), while OCLC Research found only 48% compliance among library repositories (Ayers et al., 2009).

For some time OAI-PMH enjoyed attention from Internet search engines, and Google, Yahoo and MSN all used it to harvest metadata from repositories, with varying degrees of success (McCown, Liu, Nelson, Zubair, & Liu, 2006). But because repository compliance was voluntary and inconsistent, and harvesting was slow, search engine companies soon grew disenchanted. In 2008 Google quietly announced that it would no longer support OAI-PMH, citing “information we gain from our support of OAI-PMH is disproportional to the amount of resources required to support it (Mueller, 2008). The announcement went unnoticed for some time by repository managers who had come to rely on it as a means to get their collections harvested and indexed by search engines.

Ebooks

Ask any librarian their preferred format for reading an ebook, and the answer is likely to be Amazon’s Kindle, Apple’s iBooks, or Barnes & Noble’s Nook. It’s not surprising because those choices reflect the preferences of the general population.

While the digital publishing market is in constant flux, the Book Industry Study Group (BISG) reports that Amazon holds 67% market share, Barnes & Noble 11.8% and Apple iBooks 8.2%, with the rest probably consisting of “Kobo, Google, Sony for retail, direct from publishers or perhaps the public library” (Perry, 2013). Another report states “e-books now make up around 30% of all book sales, and Amazon has a 65% share within that category, with Apple and Barnes & Noble accounting for most of the balance” (Bercovici, 2014).

“Most of the balance.” Consider that phrase for a moment. The market share of the second and third most popular ebook platforms, iBooks and Nook, are each thought to hover between 8%-20%, leaving a tiny, tiny slice for other ebook reader platforms (Wahba, 2014). Yet, these “other” platforms, which we ourselves are loath to use, tend to be the options we offer our customers.

Even the electronic version of the book my co-author and I published with a library publisher is only available for download as a PDF file. The print version appears in Amazon, but there is no Kindle version available for sale. The publisher has simply been unable to deliver to the most popular ebook reader platform in the world.

Libraries have been trying to position themselves as publishers, but it’s difficult to imagine a successful publisher who publishes to platforms where people aren’t. A new ebook lending process and reader called Occam’s Reader is currently under development by Texas Tech University, University of Hawai’I at Manoa, and the Greater Western Library Alliance in agreement with Springer. While the process may solve some longstanding problems related to lending ebooks, the reader itself doesn’t sound like an improvement on existing products: “The final reader that patron interacts with is very basic ... simple image files, no OCR, no search, only basic navigation...” (Ketner, 2013)

Implications for Library Administrators

Library administrators must constantly make decisions about resources and impact, and given the amount of labor that our data interchange practices require it is legitimate to ask about the return on investment (ROI) of those practices. While some administrators have raised justifiable concerns about libraries responding to questions of quantifying ROI (Neal, 2011), it is equally justifiable to be concerned about practices that result in relatively small audiences. Publicly funded libraries may be under more pressure to demonstrate effective use of funds, but even private institutions must answer to donors, many of whom are increasingly attuned to business acumen and attention to ROI.

Librarians and archivists may protest that crosswalks exist that allow us to work across schemas. But crosswalks are like duct tape on a problem that we know needs a more serious fix. Crosswalks always lose granularity and raise a whole new round of decisions that have to be made about mappings.

Conclusion

The thesis of this article has been that locally developed library content is not well represented on the Internet, which is the information ecosystem where most people gather their information and use it. If we believe there's value to making our materials discoverable and usable to a wider audience of people, then we must begin a concerted effort to make our metadata interoperable with Web standards and to publish to platforms that more people use. Well-intentioned librarians and archivists have spent decades developing and implementing data interchange formats that simply haven't been adopted by the Internet, and as a result we struggle to make our materials visible and usable. Is it appropriate, or even responsible, for us to continue to push an environment where people aren't, or where they don't even seem to want to be?

Henriette Avram's 1975 summary of the decades that preceded full adoption of MARC resonate today: "Libraries have passed through an era of much talk but few results into the 1970's where the automation of library operations is no longer a promise but a demonstrated success" (Avram, 1975). Today we face another transition as revolutionary as the move from the card catalog to online catalog. In a future article I hope to discuss some of the efforts that are being made to migrate library data interchange standards to interoperability with Internet search engines and other platforms.

Arlitsch, K., & O'Brien, P. S. (2012). Invisible institutional repositories: Addressing the low indexing ratios of IRs in Google Scholar. *Library Hi Tech*, 30(1), 60–81. doi:10.1108/07378831211213210

Arlitsch, K., & O'Brien, P. S. (2013). *Improving the visibility and use of digital repositories through SEO*. Chicago: ALA TechSource, an imprint of the American Library Association. Retrieved from <http://search.ebscohost.com/login.aspx?direct=true&scope=site&db=nlebk&db=nlabk&AN=578551>

Avram, H. D. (1975). MARC: its history and implications. Library of Congress. Retrieved from <http://files.eric.ed.gov/fulltext/ED127954.pdf>

- Ayers, L., Camden, B. P., German, L., Johnson, P., Miller, C., & Smith-Yoshimura, K. (2009). *What we've learned from the RLG partners metadata creation workflows survey* (No. 1-55653-409-4, 978-1-55653-409-6). Dublin, Ohio: OCLC Research. Retrieved from <http://www.oclc.org/programs/publications/reports/2009-04.pdf>
- Barnard, D. T., Burnard, L., DeRose, S. J., Durand, D. G., & Sperberg-McQueen, C. M. (1995). Lessons for the World Wide Web from the Text Encoding Initiative. In *The Web Revolution*. Boston: World Wide Web Consortium. Retrieved from <http://www.w3.org/Conferences/WWW4/Papers/337/>
- Bercovici, J. (2014, February 10). Amazon vs. book publishers, by the numbers. *Forbes*. Retrieved from <http://www.forbes.com/sites/jeffbercovici/2014/02/10/amazon-vs-book-publishers-by-the-numbers/>
- Cutts, M. (2009, September 21). Google does not use the keywords meta tag in web ranking. Retrieved from <http://googlewebmastercentral.blogspot.in/2009/09/google-does-not-use-keywords-meta-tag.html>
- Dempsey, L. (2010, January 11). Outside-in and inside-out. Retrieved from <http://orweblog.oclc.org/archives/002047.html>
- Dempsey, L., Malpas, C., & Lavoie, B. (2014). Collection Directions: The Evolution of Library Collections and Collecting. *Portal: Libraries and the Academy*, 14(3), 393–423. doi:10.1353/pla.2014.0013

- Gartner, R. (2014). An XML schema for enhancing the semantic interoperability of archival description. *Archival Science*. doi:10.1007/s10502-014-9225-1
- Google Scholar. (2011). Inclusion Guidelines for Webmasters [Inclusion Guidelines]. Retrieved October 4, 2011, from <http://scholar.google.com/intl/en/scholar/inclusion.html>
- Group, N. A. M. U., Systems, U. A. for O., Institute, E. P. R., Forum, N. A. I. U., & International, C. for O. S. (1991). *The User's Open Systems Conference, November 18-22, 1991, Reston Hyatt Regency* (Vols. 1-2). Reston, VA: Corporation for Open Systems International. Retrieved from <http://books.google.com/books?id=ULlruAAACAAJ>
- Ketner, K. (2013). *Occam's Reader*. Code4Lib. Retrieved from <https://archive.org/details/Code4libKennyKetner>
- Lagoze, C., & Van de Sompel, H. (2003). The making of the Open Archives Initiative Protocol for Metadata Harvesting. *Library Hi Tech*, 21(2), 118–128. doi:10.1108/07378830310479776
- McCown, F., Liu, X., Nelson, M. L., Zubair, M., & Liu, X. (2006). Search engine coverage of the OAI-PMH corpus. *IEEE Internet Computing*, 10(2), 66–73. doi:10.1109/MIC.2006.41
- Mueller, J. (2008, April 23). Retiring support for OAI-PMH in Sitemaps [Blog]. Retrieved from <http://googlewebmastercentral.blogspot.com/2008/04/retiring-support-for-oai-pmh-in.html>

- Mylonas, E., & Renear, A. (1999). The Text Encoding Initiative at 10: not just an interchange format anymore - but a new research community. *Computers and the Humanities*, 33(1-2), 1-9.
- National Institutes of Health. (n.d.). NIH Big Data to Knowledge (BD2K). Retrieved September 5, 2014, from <http://bd2k.nih.gov/#sthash.8px85qqq.dpbs>
- Neal, J. G. (2011). Stop the madness: the insanity of ROI and the need for new qualitative measures of academic library success. In *ACRL 2011*. Philadelphia: Association of College and Research Libraries. Retrieved from http://www.ala.org/acrl/sites/ala.org.acrl/files/content/conferences/confsandpreconfs/national/2011/papers/stop_the_madness.pdf
- OCLC, Inc. (n.d.). Data strategy and linked data. Retrieved August 20, 2014, from <http://www.oclc.org/data.en.html>
- Perry, J. W. (2013, November 14). BISG report - a few more ebook stats. Retrieved from <http://www.digitalbookworld.com/2013/bisg-report-a-few-more-ebook-stats/>
- Pitti, D. V. (1999). Encoded archival description: An introduction and overview. *New Review of Information Networking*, 5(1), 61-69.
doi:10.1080/13614579909516936
- Seeley, C. E., Nicewander, D., Page, R., & Dysert, P. A. (2004). A baseline study of medication error rates at Baylor University Medical Center in preparation for implementation of a computerized physician order entry system. *Proceedings (Baylor University Medical Center)*, 17(3), 357-361.

- Segall, L. (2014, June 12). Uber CEO: "Our growth is unprecedented." *CNN Money*.
New York: CNN. Retrieved from
<http://money.cnn.com/2014/06/12/technology/innovation/uber-ceo-travis-kalanick/>
- Styles, R., Ayers, D., & Shabir, N. (2008). Semantic MARC, MARC21 and the Semantic Web. In *Linked Data on the Web Workshops*. Beijing. Retrieved from
<http://events.linkedata.org/ldow2008/papers/02-styles-ayers-semantic-marc.pdf>
- Tibbo, H. R., & Meho, Lokman I. (2001). Finding finding aids on the world wide web. *American Archivist*, 64(1), 61–77.
- Urban, R. (2011, May 27). Beyond OAI-PMH. Retrieved from
<http://www.inherentvice.net/?p=383>
- Wahba, P. (2014, August 20). Barnes & Noble's path to e-book renaissance goes through Samsung. *Fortune*. Retrieved from <http://fortune.com/tag/nook-tablet/>
- Washburn, B., Eckert, E., & Proffitt, M. (2014, July 31). About ArchiveGrid. Retrieved August 29, 2014, from <http://beta.worldcat.org/archivegrid/about/>
- Weibel, S. L. (2010). Dublin Core Metadata Initiative: a personal history. In *Encyclopedia of Library and Information Science* (Third., Vol. 3, pp. 1655–1663). Boca Raton: CRC Press. Retrieved from
<http://www.oclc.org/research/publications/library/2009/weibel-elis.pdf>.