**DC**PAPERS

*Proc. Int'l Conf. on Dublin Core and Metadata Applications 2011*

# Describing Collections & Collection Services for the BTP

| Timothy W. Cole | Myung-Ja Han | Doug Moncur | Harriett E. Green |
|---|---|---|---|
| University of Illinois, USA | University of Illinois, USA | The Australian National University, Australia | University of Illinois, USA |
| t-cole3@illinois.edu | mhan3@illinois.edu | doug.moncur@anu.edu.au | green19@illinois.edu |

## Abstract

Libraries facilitate the use of information by collating related items to create distinct, cohesive collections. As libraries have acquired and generated more digital content, the need to agree on a standard method for describing digital collections has become increasingly evident. Shared rules for collection description not only facilitate discovery; they also have the potential to facilitate the reuse of collections and collection items. In the last decade, work has focused largely on standards and practices that facilitate collection discovery and provide human-readable descriptions of collections. With the advent of projects such as the Australian National Data Service (ANDS) and now the Bamboo Technology Project (BTP), there is a need to consider computer-mediated collection interoperability and computer-agent collection use as well. This requires more attention in collection descriptions to machine-actionable descriptions of collection-level services and suggests benefits possible through greater reliance on Semantic Web technologies such as the Resource Description Format (RDF). Experience from the Institute of Museum of Library Services Digital Collections and Content (IMLS DCC) project at University of Illinois also indicates that content-providers on their own typically do not produce collection-level descriptions that are adequate for some functions that aggregators want to deploy. This suggests that the creation of collection-level descriptions should be a collaborative enterprise. In the context of the BTP, this paper discusses the current practice of creating collection-level descriptions and introduces new developments and emerging approaches which can drive and support collection content interoperability at a more robust level.

**Keywords:** collection description; collection description application profiles; collection-level service description; ANDS; RIF-CS; Bamboo Technology Project; RDF; digital humanities

## 1. Background & Current Practice

The importance of describing digital collections in aggregated environments comes up often in the recent literature. Hill and Janee (1999) describe collection-level descriptions as a pre-requisite for some digital library services. Heaney (2000) suggests that collection descriptions help users to see and navigate the information landscape. Others have described ways that collection-level descriptions can enhance resource discovery and user satisfaction in environments where content or metadata has been aggregated across heterogeneous collections or in-cross domain contexts (e.g., Foulonneau et al., 2005; Cole & Shreeves, 2004). Chapman (2004) found that registry services based on collection descriptions can improve management of metadata in aggregation-based service environments. Brockman, et al. (2001) argue that ease of access will consistently be a factor in scholars' choice of materials, and digital libraries for scholarship require services that assist in the development and federation of collections. Another persistent theme throughout the literature argues the importance of understanding user practices as a way to inform the structure and functionality of digital collections. Collection description should reflect the way scholars use digital resources and also should adapt to the evolving needs of scholars and the services that support scholarly research and pedagogy. Palmer, Zavalina, and Fenlon (2010) argue cogently that mass aggregations of items are not enough; it is critical to retain the context and identity of collections within repositories and aggregations. Davenport (2007) asserts that digital collections useful for scholarship require good collection descriptions

**DC**PAPERS

*Proc. Int'l Conf. on Dublin Core and Metadata Applications 2011*

as a complement to expert curation of content. Norica (2007) found that in order for digital libraries to be used in teaching and learning, collection-level descriptions developed with scholarly expertise are required to enable scholars to more easily navigate resources. Maron, Smith and Loy (2009) said that a critical best practice for digital library sustainability is creating tools that provide added value to users. Sustainability of digital libraries relies on the maintenance of collection descriptions that add value for users.

So, it is not surprising then that both content providers and aggregators have been engaged in developing application profiles for collection-level descriptions. Projects looking at shared services also have suggested profiles that touch on collection-level description and the description of collection-level services. Heaney (2005) provides the generic motivation for describing collection-level services, suggesting the usefulness of such services in helping the end-user. While traditional metadata formats such as EAD and MARC allow for description of some collection attributes, these formats were not created primarily for this purpose, nor optimize to provide collection-level service description. For this paper, we identified seven collection description application profiles as representative of relevant work to date. (Links for the seven collection-level description application profiles can be found in References.)

When looking at these application profiles together, three observations emerge: first, the attributes currently seen as core are similar across profiles and focus on facilitating the discovery of collections. Second, most profiles focus on human-readable descriptions of collections, i.e., include few machine-readable attributes. And lastly, with the exception of ANDS Registry Interchange Format - Collections and Services (RIF-CS) and to a lesser extent Joint Information Systems Committee Information Environment Service Registry (JISC IESR) and Ockham, most collection-level description schemes that do provide semantics for mentioning service URLs do not type these URLs, or include sufficiently granular attributes for defining how to use these URLs. Table 1 summarizes attributes used to describe collection-level services from the profiles examined.

TABLE 1: Elements used for service-level description

| Application profile | Service-level description element(s) |
| --- | --- |
| DC Collections Application Profile (DCMI) | <isAccessedVia> |
| IMLS DCC | <isAccessedVia> |
| RLSP | not available |
| IESR (JISC) | 20 elements for service information |
| RIF-CS (ANDS) | 13 elements that focused on discovery (& access) services |
| Ockham | relies on Z39.50 to provide service information |
| Z39.91 | <CollectionService> |

## 2. Evolving role of collection description in digital libraries

The motivation for aggregators to begin using collection descriptions to facilitate service delivery can be found in use cases from many different scholarly domains. Over the last decade, the social science domain has seen the emergence of large, actively curated data warehouses, such as the UK Data Archive and the Australian Social Sciences Data Archive. The impetus towards this came from the use and reuse of large-scale survey data, often sourced from health and governmental data providers. Data sets represent a special type of collection, and sharing data sets enables reuse, obviating the need for redundant data collection. However, discoverability remains as a significant problem. One data set looks much like another and the structure of the data is often opaque. Formal collection descriptions for data sets are required to differentiate one from another. Descriptive attributes must allow users to select and identify data sets based on topical coverage as well as methodology, scope of data collection, suitability for reuse in another context.

**DC**PAPERS

*Proc. Int'l Conf. on Dublin Core and Metadata Applications 2011*

Scientific disciplines increasingly faced similar pressures to create and share data sets. Because of the increasing cost of large experiments in the sciences, the idea of making data available for reuse has gained support and led to the development of specialized collections such as the International Virtual Observatory Alliance, a federation of virtual observatories in the US, UK, Europe and Australia, and the IRIS data management center in the earth sciences. In disciplines such as genomics and phenomics, scholars increasingly make data available for reuse. Like social science data sets, scientific data sets are difficult to identify and reuse without contextual information. However, whereas social science data sets are most often assembled by a close-knit team of specialists or archivists from at most a few institutions, data sets in science disciplines are built over time by loose confederations of individuals drawn from many institutions. This highlights early the importance of shared models of description and collections interoperability. Within some disciplines, this was achieved in part by reliance on standard data formats, such as the FITS image format in astronomy and the SEED format in seismology, which mandate the metadata header to describe the provenance and context of the data. However, not all data can be represented in a standard format, and rich descriptions of data sets remain a priority.

The humanities are often regarded as computationally averse and without a tradition of structuring data. However, this is no longer the case. Textual analysis, as opposed to textual criticism, is fundamentally about counting things, a task in which computers excel, and several data mining tools for texts have been developed. Other literary research projects exploit databases to connect texts and page images, e.g., Proust's *À la recherche du temps perdu* (*In Search of Lost Time)*. There has also been intriguing work in musicology, including the use of high bandwidth links to co-ordinate geographically dispersed performers. Archaeologists have used grid computing techniques by archaeologists to simulate the battle of Manzikert. Likewise, historians use large data bases such as the Anglo-Saxon Charters database to analyze changes in property ownership over time and to track the movements of individuals over time as an aid to historical analysis. In other humanities disciplines, more rich analogs to scientific data sets have been created, such as the computational analysis of the social graph of Medieval Languedoc, and the Australian Pulp Fiction data collection digitized more than 5,000 Australian pulp fiction items published between 1939 and 1959. As other disciplines, rich, machine-actionable collection descriptions are essential to insure the broad reuse of humanities data set collections.

Traditional digital library collections, e.g., cultural heritage collections such as indexed in the IMLS DCC collection registry, also benefit from high quality, complete collection-level description. Transaction log analysis of the IMLS DCC portal (http://imlsdcc.grainger.uiuc.edu/), which allows searching for collections alongside items, suggests that users make use of collection links when discovered at almost the same rate as links to items discovered. Given the importance of collection-level descriptions, IMLS DCC has moved from relying almost exclusively on content providers to specify collection-level description attributes, to a model in which the aggregator takes the lead in creating collection-level descriptions (Palmer, Zavalina, and Fenlon, 2010).

## 3. Implications for the Bamboo Technology Project

Given the project's focus on building applications and shared infrastructure for humanities research, collections interoperability is a high priority objective for the BTP. Increasingly humanities scholars have the need to apply tools across corpora that span multiple collections and multiple repositories. The recognition of similar needs across almost all disciplines led ANDS to adopt the RIF-CS as its native metadata format. RIF-CS is based on the draft ISO 2146 *Registry Services for Libraries and Related Organisations* standard and provides semantics for describing collections and associated services, activities, and parties (agents). It includes its own controlled vocabularies, including for classes of common services (e.g., *syndicate-atom*, *harvest-oaipmh*, *search-http*). A number of its core collection description element names map directly to Dublin Core elements, though RIF-CS includes many additional concepts and allows refinement and

**DC** PAPERS

*Proc. Int'l Conf. on Dublin Core and Metadata Applications 2011*

elaboration through the use of attributes and child elements. RIF-CS is rich enough to support flexible presentation of records. Figure 1 shows an ANDS collection record for a social science
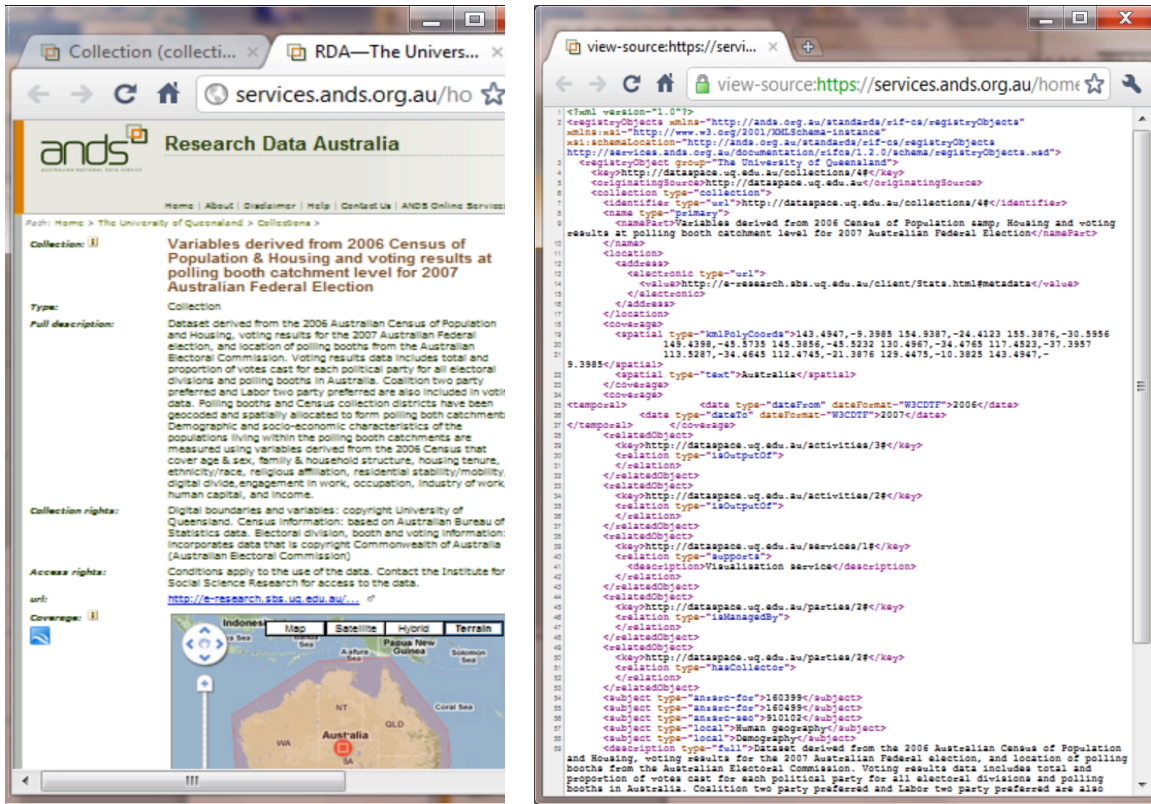


FIG. 1. Collection Record as displayed in *Research Data Australia* (left) and as an xml instance (right)
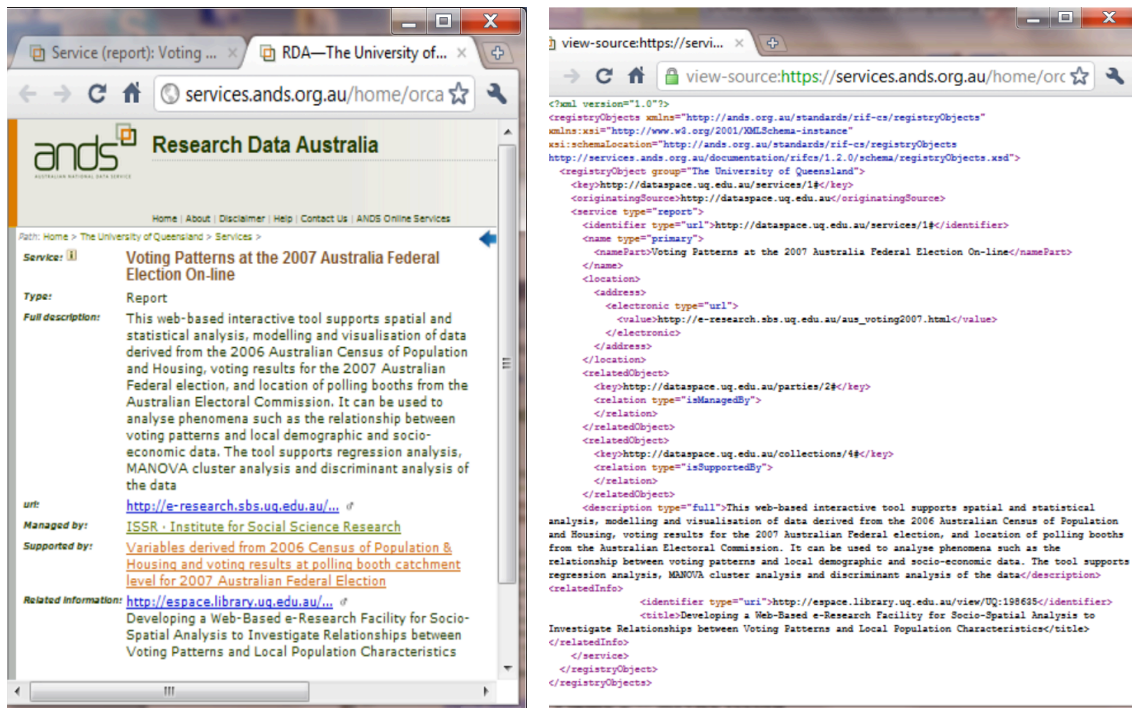


FIG.2. RIF-CS service record in user interface (left) and as an xml instance (right)

census data set as displayed in *Research Data Australia* (http://services.ands.org.au/pages/) (left) and as an XML instance (right). Note the actionable spatial coverage information included in the record. Figure 2 shows the RIF-CS record for a collection-level service associated with this data set (left) and as an xml instance (right). As illustrated in the left-hand portion of Figure 3, collection service records specify service type and give information about arguments, both required and optional.

To meet the needs of the BTP, we anticipate needing to extend the RIF-CS semantics and controlled vocabulary to describe classes of more specialized services of interest specifically to researchers and tools, e.g., a CMIS-compliant service for item dissemination, a service to deliver morphologically adorned texts from the collection, etc. Additionally, RIF-CS does not conform to the RDF data model; to best support BTP collections interoperability goals and be consistent with other humanities-centric projects like 18thConnect (http://www.18thconnect.org/), we see this as a limitation. Fortunately, though not natively RDF-compliant, we are optimistic that it will be possible to leverage the semantics and underlying data model of RIF-CS in an RDF-compliant version not far removed from the spirit and intent embodied in the current RIF-CS schema. Figure 3 shows an example of an RIF-CS service description record transformed to be consistent with the RDF data model. However, we also recognize that extending RIF-CS vocabularies and transforming the scheme to be more RDF friendly will further complicate what is already a relatively complex format for describing collections and collection-related services. Accordingly we see our use of RIF-CS as primarily limited to the back-end of Bamboo infrastructure.

```
<service type="syndicate-rss">              <rif-cs:Service>
  <name type="primary">                       <rif-cs:serviceType>syndicate-rss</rif-cs:serviceType>
    <namePart>                                <rif-cs:hasName>
      RSS 2.0 Feed from                         <rif-cs:Name>
      MY University institutional Repository      <rif-cs:nameType>Primary</rif-cs:nameType>
    </namePart>                                   <rif-cs:namePart>
  </name>                                           RSS 2.0 Feed from MY University Institutional Repository
  <location>                                      </rif-cs:namePart>
    <address>                                   </rif-cs:Name>
      <electronic type="url">                 </rif-cs:hasName>
        <value>http://myrepo.myu...           <rif-cs:hasLocation>
        </value>                                <rif-cs:ElectronicAddress rdf:about="http://myrepo.myu....">
        <arg required="true"                      <rif-cs:electronicAddressType>url
        type="string" use="keyValue">           </rif-cs:electronicAddressType>
          identifier</arg>                      <rif-cs:hasArg>
      </electronic>                               <rif-cs:Arg>
    </address>                                     <rif-cs:required>true</rif-cs:required>
  </location>                                      <rif-cs:argType>string</rif-cs:argType>
</service>                                         <rif-cs:argUse>keyValue</rif-cs:argUse>
                                                  <rif-cs:argName>identifier</rif-cs:argName>
                                                </rif-cs:Arg>
                                              </rif-cs:hasArg>
                                            </rif-cs:ElectronicAddress>
                                          </rif-cs:hasLocation>
                                        </rif-cs:Service>
```

FIG. 3. Collection-level description about collection service in RIF-CS (left) and RDF (right)

Finally, for the BTP we anticipate that the creation of collection-level description takes place over time and involves multiple authors. As metadata sharing and aggregation has become the norm in the digital library setting, service providers have also begun to create collection-level descriptions for the collections they aggregate. For the BTP, collection-level description will likely be gathered initially from content providers using more ubiquitous schemas, e.g., the Dublin Core Collection Description Application profile, then converted through automatic means to an extended version of RIF-CS and augmented by human cataloger intervention.

## DCPAPERS

*Proc. Int'l Conf. on Dublin Core and Metadata Applications 2011*

## References

ANDSCMS. (n.d.). Research Data Australia. Retrieved July 17, 2011, from http://www.ands.org.au/guides/research-data-australia.html

Australian National Data Service. (n.d.). RIF-CS Guide. Retrieved from http://ands.org.au/guides/rif-cs-awareness.html

Brockman, William S., Laura Neumann, Carole L. Palmer, and Tonyia J. Tidline. (2001). Scholarly Work in the Humanities and the Evolving Information Environment. Retrieved July 17, 2011, from http://www.clir.org/ pubs/ reports/pub104/pub104.pdf

Chapman, Ann. (2004). Collection-level description: joining up the domains. *Journal of the Society of Archivists*, *25*(2), 149-155. doi:10.1080/0037981042000271475

Cole, Timothy W. and Sarah L. Shreeves. (2004). Search and discovery across collections: the IMLS digital collections and content project. *Library Hi Tech*, 22(3), 307-322. doi:10.1108/07378830410560107

Davenport, Nancy. (2007). Digital Libraries and Librarians of the 21st Century. *Journal of Library Administration*, 46(1), 89-97. doi:10.1300/J111v46n01-07

Dublin Core Collection Description Working Group. (2006). Dublin Core collection description application profile. Retrieved July 17, 2011, from http://dublincore.org/groups/collections/collection-application-profile/

Foulonneau, Muriel, Timothy W. Cole, Thomas G. Habing, and Sarah L. Shreeves. (2005). Using collection descriptions to enhance an aggregation of harvested item-level metadata. Proceedings of the ACM/IEEE Joint Conference on Digital Libraries, 32-41.

Heaney, Michael. (2000). An analytical model of collections and their catalogues. Retrieved July 17, 2011, from http://www.ukoln.ac.uk/metadata/rslp/model/amcc-v31.pdf

Heaney, Michael. (2005). Users and information resources: an extension of the analytical model of collections and their catalogues into usage and transactions. Retrieved July 17, 2011, from http://www.ukoln.ac.uk/cd-focus/model-ext/CD2-principles-v2-2.pdf

Hill, Linda L., and Greg Janee. (1999). Collection metadata solutions for digital library applications. Journal of the American Society for Information Science, 50(13), 1169-1181.

Joint Information Systems Committee. (2010). Information environment service registry (IESR). Retrieved July 17, 2011, from http://www.jisc.ac.uk/whatwedo/services/mimas/iesr.aspx

Maron, Nancy L., Kirby Smith, and Matthew Loy. (2009). *Sustaining digital resources:an on-the-ground view of projects today: Ithaka case studies in sustainability.* Retrieved July 17, 2011, from http://www.ithaka.org/ithaka-s-r/strategy/ithaka-case-studies-in-sustainability

National Information Standards Organization. (2005). Z39.91 - 200X - Collection Description Specification. Retrieved July 17, 2011, from http://www.niso.org/kst/reports/standards/kfile_download?id%3Austring%3Aiso-8859-1=Z39-91-DSFTU.pdf&pt=RkGKiXzW643YeUaYUqZ1BFwDhIG4-24RJbcZBWg8uE4vWdpZsJDs4RjLz0t90_d5_ymGsj_IKVa86hjP37r_hAkGN_NSjktlVeUOy23hQur71o0A7vDdEIoCbA1-kWBRkh0FbjkOU7U%3D

Norcia, Megan. A. (2008). Out of the ivory tower endlessly rocking: collaborating across disciplines and professions to promote student learning in the digital archive. *Pedagogy: Critical Approaches to Teaching Literature, Language, Composition, and Culture*, 8(1), 91-114. doi:10.1215/15314200-2007-026

OCKHAM Overview. (2006). Retrieved July 17, 2011, from http://ockham.org/

Palmer, Carole, Oksana Zavalina, and Katrina Fenlon. (2010). Contextual mass in digital aggregations for scholarly use. *Proceedings of the American Society for Information Science and Technology* 47 (1): 1-10. DOI: 10.1002 /meet.14504701213

Powell, Andy, Michael Heaney, and Lorcan Dempsey. (2000). RSLP collection description. *D-Lib Magazine*, 6(9). Retrieved, July 17, 2011, from http://www.dlib.org/dlib/september00/powell/09powell.html

University of Illinois IMLS Digital Collections and Content Project. (2008). Retrieved July 17, 2011, from http://imlsdcc.grainger.uiuc.edu/cdschema_overview.asp