

## Describing OAI Sets: A Discussion Paper

### ► Introduction to issues around OAI set descriptions

The Digital Library Federation and the National Science Digital Library gathered in July 2004 to lead a joint initiative to review current issues with OAI implementations and define best practices for OAI and sharable metadata.<sup>1</sup> The DLF / NSDL best practices for OAI and sharable metadata group agrees that:

- Set and collection descriptions are not the same although there may be overlap between the two;
- Service providers find set descriptions useful;
- Service providers also find collection descriptions useful as distinct from set descriptions;
- The OAI implementation guidelines acknowledge that there is a difference between set and collection descriptions but do not explicitly explain the differences. The implementation guidelines refer to both collection description formats, including EAD and the Dublin Core Collection Description Application Profile, and service description formats such as the ePrints schema and UDDI when talking about set description<sup>2</sup>;
- Data providers currently provide many different types of information within the set description. Some of this information describes collections of resources and some of which describes the OAI set as distinct from a collection of resources;
- In general, data providers do not want to repeat the same information for sets and collections;

Therefore, we believe there needs to be:

- 1) a clarification of the difference between a description of an OAI set and a description of a collection of resources;
- 2) a metadata format or an application profile to specifically describe an OAI set;
- 3) a recommendation on the mechanism to convey both the set description and collection description.

This discussion paper proposes a framework for both and seeks comments about the content presented here. We are actively seeking comments and feedback and welcome anyone who might want to participate in the further development of this proposal. Please contact Muriel Foulonneau ([mfoulonn@uiuc.edu](mailto:mfoulonn@uiuc.edu)) and Sarah Shreeves ([sshreeve@uiuc.edu](mailto:sshreeve@uiuc.edu)) with comments, questions, and feedback.

---

<sup>1</sup> [http://oai-best.comm.nsdsl.org/cgi-bin/wiki.pl?OAI\\_Best\\_Practices](http://oai-best.comm.nsdsl.org/cgi-bin/wiki.pl?OAI_Best_Practices)

<sup>2</sup> <http://www.openarchives.org/OAI/2.0/guidelines-repository.htm#setDescription>

► **What is a collection?**

A collection is “any aggregation of physical or digital items”<sup>3</sup>. Two types of collections are discussed here: the collection of metadata (an OAI set) and a collection of resources.

► **What is an OAI set?**

Within the OAI environment, “a set is optional construct for grouping items for the purpose of selective harvesting.”<sup>4</sup> An OAI set is both:

- a collection of metadata records in an OAI repository; and
- the technical mechanism by which the OAI repository (the service) makes these metadata records available.

An OAI set may include a set description. **We propose that a set description should describe both the collection of metadata records as well as the appropriate elements of the technical mechanism.**

► **Elements of a set description**

**We propose to use a Dublin Core application profile to describe OAI sets.**

The list below was created from the list of elements that are in actual use (study from the UIUC OAI registry<sup>5</sup>) and are also drawn from the last version of the Dublin Core Collection Description Application Profile (CD AP).<sup>6</sup> In addition, the OAI Best Practices Working Group has considered the elements that would be useful for service providers for the purposes of selective harvesting and processing of metadata records.

The sections in italics and red are areas where we are particularly interested in comments from the community.

(Note that not all elements indicated below have been officially accepted into the Dublin Core Collection Description Application Profile; the cld and gen namespaces used below are temporary until the CD AP is formalized.)

Description	Potential Metadata Element	Example	Comment
A unique identifier in the OAI world	dc:identifier	- lcoa1.loc.gov:brhc	This is the repository identifier + oai setspec
Set name	dc:title	- Brady-Handy Collection (Photographs)	This is a repetition of the oai:setname

<sup>3</sup> <http://www.dublincore.org/groups/collections/>

<sup>4</sup> <http://www.openarchives.org/OAI/openarchivesprotocol.html#Set>

<sup>5</sup> <http://gita.grainger.uiuc.edu/registry/>

<sup>6</sup> <http://www.ukoln.ac.uk/metadata/dcmi/collection-ap-summary/>

Metadata formats in which the set is available	- dc:format OR - cld:itemFormat	- MODS - oai_dc	Please note that the use of dc:format within the DC Collection Description AP is under discussion. See the DC CD listserv in May-July 2005 for further discussion. <sup>7</sup>
How records are added to the set	dct:accrualMethod	- Updated automatically from a database - Updated through manual export	
Whether the set is updated	dct:accrualPolicy	- No longer adding items - Closed - Records are modified	This should be coherent with the update frequency.
Update frequency	dct:accrualPeriodicity	- Monthly - Weekly	Allows service providers to know how often to reharvest.
Earliest and latest datestamps within the set	cld:contentDateRange	- 2004-11-03/2005-03-05	Allows a service provider to know whether to reharvest.
Size of the set	dct:extent	- 1000 records	Allows a service provider to double check number of harvested records
Person in charge of the metadata records / responsibility	marcrel:own	- John Smith <jsmith@collectionsinstitute.org>	This might be useful if the institution that created the records is in a different institution than the institution that hosts the OAI data provider.
Person in charge of the set maintenance / responsibility	dc:creator	- Elizabeth Taylor <etaylor@oaiservices.net>	Presumably this should be the same as the administrator in the OAI Identify response.
Access rights	dct:accessrights	- HTTP authentication required, by prior agreement – contact Elizabeth	Could be used if an authentication is required to harvest that specific set.

<sup>7</sup> <http://www.jiscmail.ac.uk/lists/DC-COLLECTIONS.html>

		Taylor <etaylor@oaiservices.net>	
<i>Rights</i>	<ul style="list-style-type: none"> <li>- <i>dc:rights</i></li> <li>OR</li> <li>- <i>&lt;rightsManifest&gt;</i></li> </ul>	<ul style="list-style-type: none"> <li>- <i>See &lt;rights&gt; container for rights held over the metadata in this set.</i></li> <li>- <i>&lt;rightsManifest&gt;...&lt;/rightsManifest&gt;</i></li> </ul>	<i>Within the OAI context, rights for the metadata records (NOT resources) are defined in the &lt;rights&gt; container. The rights element should refer users to this container or could alternatively duplicate the rights container. Please note that there is no way to express rights over the set (aggregation) of metadata records; only the individual records.</i>
<i>A reference to a collection description of the resources described by the metadata, if applicable</i>	<i>dct:references</i>	<ul style="list-style-type: none"> <li>- <i>http://imlsgcc.grainger.uiuc.edu/collections/FullDisplay.asp?cid=2391</i></li> <li>- <i>&lt;cid:collection&gt;...&lt;/cid:collection&gt;</i></li> </ul>	<i>This is where a reference to a collection description should be contained through any number of ways (see below for a further discussion). URLs to collection descriptions? oai:identifier of collection description records in the repository? Embedding collection descriptions?</i>
<i>Summary statement about the content of the set</i>	<i>dct:abstract</i>	<ul style="list-style-type: none"> <li>- <i>Records for cookbooks digitized as part of the Feeding America: Historic American Cookbooks Project</i></li> </ul>	
<i>Documentation about the metadata records including cataloging and mapping info</i>	<i>dc:description</i>	<ul style="list-style-type: none"> <li>- <i>mapped to MODS from MARC records</i></li> <li>- <i>LCSH terminology used in MODS format</i></li> </ul>	<i>May also include a URI pointing to documentation.</i>

Relationship to other sets	<ul style="list-style-type: none"> <li>- dct:isPartOf (supercollection)</li> <li>- dct:hasPart (subcollection)</li> <li>- dc:relation (associated collection)</li> </ul>	<ul style="list-style-type: none"> <li>- Overlaps with set XXX and YYY</li> <li>- Belongs to superset DDD</li> <li>- Is composed notably although not exclusively of subsets RRR and UUU</li> </ul>	
Language of records in set	dc:language	<ul style="list-style-type: none"> <li>- en</li> <li>- fr</li> <li>- ru</li> <li>- ge</li> </ul>	
Audience	dct:audience	<ul style="list-style-type: none"> <li>- Created for the National Science Digital Library &lt;http://nsdl.org/&gt;</li> </ul>	Note specific service providers sets have been created for.
Logo	<ul style="list-style-type: none"> <li>- <i>cld:logo</i></li> <li>OR</li> <li>- <i>&lt;branding&gt;</i></li> </ul>	<ul style="list-style-type: none"> <li>- <i>see &lt;branding&gt; container for logo for this set</i></li> <li>- <i>&lt;branding&gt;</i></li> <li>...</li> <li>- <i>&lt;/branding&gt;</i></li> </ul>	<i>Similar problem as with rights (though rather less serious!)</i>
Institution which hosts the set	gen:isLocatedAt	<ul style="list-style-type: none"> <li>- Oaiservices Inc.</li> </ul>	Could be a static repository gateway for example (also the repository though)
Service that makes the metadata records available	gen:isAccessedVia	<ul style="list-style-type: none"> <li>- http://myOAIrepository/?verb=ListRecords&amp;metadataPrefix=oai_dc&amp;set=brhc</li> </ul>	Either the baseURL OAI query into the set and alternate ways to get into a set

► **Set – collection relationships**

There is no *a priori* relationship between a collection of resources and an OAI set. However, it can be useful if there is, in fact, a relationship between an OAI set and a collection of resources for a variety of reasons, including:

- Service providers can use collection information to provide richer context for the harvested metadata;
- Information that service providers might find useful for selective harvesting purposes such as access rights may be attached to a specific collection of resources but not to an entire OAI repository; and

- The metadata for resources exposed via OAI may be managed at a collection level and thus organizing sets based on collections may ease the related management of the metadata in the OAI repository.

**We propose that when possible data providers should create sets corresponding to individual collections and optionally should provide other hierarchical or overlapping sets.** A collection of resources can be defined in any way imaginable (for example, by topic, type of material, author, provenance, etc). Regardless, data providers might think about the following when defining sets in the OAI context:

- 1) The set creates a granularity unit that differentiates it from the other groupings within the context of a specific application.

For example, the Library of Congress provides multiple collections via OAI that broadly fit under the category of American History; however, these have been usefully broken apart into individual sets that correspond to specific defined collections of resources. While one service provider might be interested in the "Ansel Adams's Photographs of Japanese-American Internment at Manzanar", it may not be interested in "Glass negatives from the Papers of Wilbur and Orville Wright".

At the opposite extreme, the OAI repository for the Illinois aerial photograph collection is divided into sets by both year and county. These could be grouped into sub-collections based on the year the photographs were taken (i.e. to create a collection of aerial photos of Illinois in 1940, in 1941, etc). In our experience, these are not different enough from one another to warrant representing these in distinct sets.<sup>8</sup>

- 2) The collection of resources was created through the same digitization project or process. This often means that the metadata were at least reviewed and processed at the same time, and the formats of resources are similar. These collections would be good units for aggregators to reprocess and manage because of their potential consistency.

Obviously, it is not possible for data providers to organize resources into collections that will be useful for all harvesters, but data providers should spend some time thinking about how their sets are organized in relation to collections.

**When a set corresponds to a collection of resources, the collection description corresponding to that collection of resources should be either embedded in or referred to in the set description in the <dc:references> element.**

If sets represent multiple collections, the set description could embed or refer to multiple collection descriptions. If sets are not designed according to a concept of

---

<sup>8</sup> <http://ciharvest.granger.uiuc.edu/documents/usingcollectiondescriptions.pdf>

collection, the set description might leave the <dc:abstract> blank or include some other sort of description such as: Contains metadata records with subject: physics.

### ► **Sharing collection descriptions in the OAI environment**

The mechanism by which data providers should share collection description is not certain. While we have not fully explored how collection descriptions should be shared within the OAI environment, several alternatives present themselves each with pros and cons:

- 1) An oai record with a collection description is included in the repository, but in a set containing only collection descriptions rather than in the same set as the records describing its component items. This is the example followed by the National Science Digital Library. Potentially the collection description could be located by a GetRecord link in the set description.
- 2) A link to a collection description on a web page. This could be embedded in the set description and/or the item level metadata records; or
- 3) A metadata record that could be embedded in the set description.

An additional solution consists of creating an oai record with a collection description contained within the same set as its [the collection's] component parts. The collection description could be located by a GetRecord link in the set description. The relationship between the item level metadata records and the collection description record could be noted in the <dc:relation> field. This solution is the least palatable from both a data and service provider perspective as it requires additional maintenance of not only a collection description record, but also item level records.

**Therefore, the present proposal considers the descriptions of collections of resources as part of the <dc:references> field of the set description and suggests using either an embedded collection description or a URI to reference the collection description. The URI can be either a Webpage <http://myWebSite/mycollectiondescription.html> or an oai record <http://myOAIrepository.org/?verb=GetRecord&metadataPrefix=dccoll&identifier=oai:MyCollectionDescriptionIdentifier>.**

We would welcome comments and thoughts on these methods.

### ► **Examples**

#### **From the Library of Congress**

Information about the set is included in a <dc:description> tag within a full collection description:

```
<dc:description>Set characteristics for calbkbib: Source records are MARC (from LC catalog); MODS or oai_dc records are dynamically generated using generic transformation when harvested.
```

dct:accrualPolicy: Closed. Contains about 200 records. Records in set calbkbib are also in set lcbooks.</dc:description>

### **From Michigan State University:**

Models the proposal in this paper:

```
<set>
  <setSpec>fap</setSpec>
  <setName>Feeding America: The Historic American Cookbook Project
  </setName>
  <setDescription>
    <dc:identifier>lib.msu.edu:fap</dc:identifier>
    <dc:title>Feeding America: The Historic American Cookbook
    Project</dc:title>
    <dc:format>oai_dc</dc:format>
    <dct:accrualmethod>dynamically generated from local
    database</dct:accrualmethod>
    <dct:accrualpolicy>Records may be added as new cookbooks
    are digitized; records may be modified if
    necessary</dct:accrualpolicy>
    <dct:accrualperiodicity>Infrequently</dct:accrualperiodicit
    y>
    <cld:contentdaterange>2005-05-20/2005-05-
    20</cld:contentdaterange>
    <dct:extent>77 records</dct:extent>
    <marcrel:own>Anne Karle-Zenith</marcrel:own>
    <dct:collector>Michael Seadle</dct:collector>
    <dct:accessrights>No restrictions</dct:accessrights>
    <dc:rights>No restrictions</dc:rights>
    <dct:references>
      <cld:collection>
        <dc:identifier>http://digital.lib.msu.edu/project
        s/cookbooks</dc:identifier>
        <dc:title>Feeding America: The Historic American
        Cookbook Project</dc:title>
        <dct:abstract>Online collection of some of the
        most important and influential American cookbooks
        from the late 18th to early 20th
        century.</dct:abstract>
        <dct:extent>77 items</dct:extent>
        <dc:language xsi:type="ISO639-
        2">eng</dc:language>
        <dc:type xsi:type="cldtype">Collection of
        Texts</dc:type>
        <dct:accessRights>No
        restrictions</dct:accessrights>
        <cld:accrualMethod
        xsi:type="DCCDAccrualMethod">ItemCreation</cld:ac
        crualMethod>
```

```
<cld:accrualPeriodicity>Unknown</cld:accrualPeriodicity>
<cld:accrualPolicy
xsi:type="DCCDAccrualPeriodicity">Passive</cld:accrualPolicy>
<dc:subject xsi:type="LCSH">Cookery,
American</dc:subject>
<dc:subject xsi:type="LCSH">Cookery -- United
States -- 19th century</dc:subject>
<dc:subject xsi:type="LCSH">Cookery -- United
States -- 20th century</dc:subject>
<dc:creator>Michigan State University Libraries.
Digital & Multimedia Center.</dc:creator>
<marcrel:own>Michigan State University Libraries.
Digital & Multimedia Center.</marcrel:own>
<gen:isLocatedAt>Michigan State University, 100
Library, East Lansing, MI, 48224,
USA</gen:isLocatedAt>
<gen:isAvailableVia
xsi:type="URI">http://digital.lib.msu.edu/projects/
cookbooks</gen:isAvailableVia>
</cld:collection>
</dct:references>
<dct:abstract >Records for cookbooks digitized as part of
the Feeding America: Historic American Cookbooks
Project</dct:abstract>
<dc:description>Simple Dublin Core records created from
documentation for each cookbook; oai_dc subject uses LCSH;
oai_dc:type uses DCMIType vocabulary; oai_dc:coverage (1)
uses TGN</dc:description>
<dc:language xsi:type="dct:ISO639-2">eng</dc:language>
<gen:isLocatedAt>Michigan State University Libraries
Digital & Multimedia Center</gen:isLocatedAt>
<gen:isAccessedVia xsi:type="dct:URI">
http://oai.lib.msu.edu/OAIHandler </gen:isAccessedVia>
</setDescription>
</set>
```