JISC CETIS Metadata and Digital Repository SIG meeting, Manchester

16 April 2007

A Dublin Core Application Profile for Scholarly Works (eprints)

Julie Allinson

Repositories Research Officer UKOLN, University of Bath





A centre of expertise in digital information management

UKO Supported by

ISC

www.ukoln.ac.uk



Museums, Libraries and Archives Council

overview

- background, scope and functional requirements
- the model
- the application profile and vocabularies
- oai-pmh , dumb-down and community acceptance



background, scope and functional requirements



terminology

- eprints, research papers and scholarly works are used synonymously for
 - a "scientific or scholarly research text"

 (as defined by the Budapest Open Access Initiative
 www.earlham.edu/~peters/fos/boaifaq.htm#liter
)
 - e.g. a peer-reviewed journal article, a preprint, a working paper, a thesis, a book chapter, a report, etc.
- the application profile is independent of any particular software application



the problem space

- simple DC is insufficient to adequately describe eprints
- the metadata produced is often inconsistent and poor quality
- identifying the full-text is problematic
- this poses problems for aggregator services



the work

- the work aimed to develop:
 - a Dublin Core application profile for eprints
 - containing properties to support functionality offered by the Intute repository search service, such as fielded searches of the metadata or indexing the full-text of the research paper;
 - any implementation / cataloguing rules;
 - a plan for early community acceptance and takeup, bearing in mind current practice
- co-ordinated by Julie Allinson (UKOLN) and Andy Powell (Eduserv Foundation), summer 2006
- through a working group and feedback group
- using a wiki to make all documentation freely available, at all times



the scope

- as provided by JISC, the funders
 - DC elements plus any additional elements necessary
 - identifiers for the eprint and fulltext(s), and related resources
 - hospitable to a variety of subject access solutions
 - additional elements required as search entry points
 - bibliographic citations and references citing other works



the functional requirements : a selection

- richer metadata set & consistent metadata
- unambiguous method of identifying full-text(s)
- version identification & most appropriate copy of a version
- identification of open access materials
- support browse based on controlled vocabularies
- OpenURL & citation analysis
- identification of the research funder and project code
- identification of the repository or service making available the copy
- date available
- date of modification of a copy, to locate the latest version



the requirements demanded a more complex model ...

the model



what is an application model?

- the application model says what things are being described
 - the set of **entities** that we want to describe
 - and the key relationships between those entities
- model vs. Model the application model and the DCMI Abstract Model are completely separate



 the DCMI Abstract Model says what the descriptions look like

FRBR

- FRBR (Functional Requirements for Bibliographic Records) is a model for the entities that *bibliographic records* are intended to describe
- FRBR models the world using 4 key entities: Work, Expression, Manifestation and Item
 - a work is a distinct intellectual or artistic creation.
 A work is an abstract entity
 - an **expression** is the intellectual or artistic realization of a work
 - a manifestation is the physical embodiment of an expression of a work
 - an **item** is a single exemplar of a manifestation.
 The entity defined as item is a concrete entity



RBR relationships

 FRBR also defines additional entities that are related to the four entities above - 'Person', 'Corporate body', 'Concept', 'Object', 'Event' and 'Place' - and relationships between them

• the key entity-relations appear to be:

- Work -- is realized through --> Expression
- Expression -- is embodied in --> Manifestation
- Manifestation -- is exemplified by --> Item
- Work -- is created by --> Person or Corporate Body
- Manifestation -- is produced by --> Person or Corporate Body
- Expression -- has a translation --> Expression
- Expression -- has a revision --> Expression
- Manifestation -- has an alternative --> Manifestation



FRBR for eprints

- FRBR provides the basis for our model
 - it's a model for the entities that *bibliographic records* describe
 - but we've applied it to scholarly works
 - and it might be applied to other *resource types*
- FRBR is a useful model for eprints because it allows us to answer questions like:
 - what is the URL of the most appropriate copy (a FRBR item) of the PDF format (a manifestation) of the pre-print version (a expression) for this eprint (the work)?
 - are these two copies related? if so, how?



the model





the model



vertical vs. horizontal relationships





an example - a conference paper

texts in whole or h issues of provenance will become being usues of the possibility of malformed day, abuse, such as spanning, unado Signed metadata: method and application This paper discusses the role that PKI digital setime that accountancy to be build 1.1 Principles behind the digital aguator The digital signature, first introduced in Diffe a The week techniques, that aims to permit the operates analogously to a handwritten sign Abstract As metadata providers increase in number and diversity, and additional contexts for metadata use an identified issues of trust, provenance and identity pain in relevance. Use of a Dublic Kee Jain an identified issues of trust, provenance and identity pain in relevance. As metadata providers increase in number and diversity, and additional contexts for metadata use are identified, issues of most provenance and identity gain in relevance. Use of a Public Key Infra-are identified, issues of most provenance and identity gain in relevance, movining evidence of the structure (PKI) is discussed for divibal simulation of metadata records, movining evidence of a structure (PKI) is discussed for divibal simulations of metadata records. operative as received is equal to the most are identified, issues of trust, procenance and identity gain in relevance. Use of a Public-Key Infra-structure (PKD) is discussed for digital signature of metadata records, providing endeneds an information within the some rand the authenticity of the information within the record. Two metadata identity of the some rand the authenticity of the information within the record, providing endeneds are mesanes signature also requires as a presquisie structure (PhD) is discussed for digital signature of metadata records, providing evidence of the identity of the signer and the authenticity of the information within the record. Two methods are suggested, firstly the W3C XML Sugnature, and secondly, identification of a minimal set of meta-suggested, firstly the W3C XML Sugnature, and secondly, identification of a minimal set of metaverify that a handwritten signame be identity of the signet and the authenticity of the information within the record. Two methods are suggested, firstly, the WZC XML Signature, and secondly, identification of a minimal second s existing signature. A similar prerousan suggested, firstly, the WSC XML sugnature, and secondly, identification of a minimal set of meta-data dements that enable signature ventication across various character sets and formats, using the OpenPGP standard. Possible strategies for bandling apponition within this infrastructure are sup-OpenPGP standard. Possible strategies for bandling apponition within this infrastructure are sup-openPGP standard. use of the digital signature is as an er data dements that enable signature verification across various character sets and formats, using the OpenPGP standard. Possible strategies for handling annotation within this infrastructure are sug-pested. Finally, some use cases are briefly discussed. relative to known identities, and to be This is possible using a variety of roday is based around PKL in wh creating the digital signature is those who need the ability to ve gested. Finally, some use cases are briefly discussed. beterogeneous infrastructure, digital signature, web of trust, provenance. public key cannot be used to si **1. Introduction** The issue of must, a level of confidence in a source, is of great importance on the Internet in mount. The source of a nince of information is a with detail in analysis is the source in mount. any member of the public to a The issue of trust, a level of confidence in a source, is of great importance on the Internet in general. The source of a piece of information is a vital detail in analysis is the provide known? Do they generally provide accurate information? Do they have a reason to provide in general. The source of a piece of information is a vital detail in analysis; is the source of a piece of information? Do they have a reason to provide accurate information? Do they generally provide accurate information? In this manner, the provenance of a piece of information provide information? In this manner, the provenance of a piece of information provide information? the private key remains secret known? Do they generally provide accurate information? Do they have a reason to provide accurate information? The provenance of a piece of information becomes inaccurate information? In this manner, the provenance of a piece of information becomes a precessary detail in analysis and interpretation. Public keys are often distri and corresponding identit The predominance of the client-server model means that this issue may often be ignored in the distribution emissionment in that mendum whole particularly in the distribution emissionment in that mendum required. It is of course pr the predominance of the client-server model means that this issue may often be ignored cither partially or wholly, particularly in the digital library environment, in that metadata providers are considered to be responsible for the accuracy of their content. effict partially or wholly, particularly in the digital library environment, in that metadua Providers are considered to be responsible for the accuracy of their content. Archives Initia Providers are considered to be responsible for the accuracy of their open Archives Initia established either implicitly, or explicitly stated within metadati, the Open Archives Initia identity information, in i Providers are considered to be responsible for the accuracy of their content. Provenance is initial initial exploring the open Archives in the open Archives in the open accuracy of their content and the open archives in the open archives in the open archives are provided to be responsible for the accuracy of their content. Provenances in the accuracy of their content, the open archives in the open archives are considered to be responsible for the accuracy of their content. Provenances in the accuracy of their content. Provenances is a set of the accuracy of their content. Provenances is a set of the accuracy of the accuracy of the open archives in the accuracy of the open archives in the accuracy of the accur a necessary detail in analysis and interpretation. logy is not a general se stablished either implicitly, or explicitly stated within metadata; the Open Archives Impa relation of metadata across systems, and a state of the DART project (Dahlquist effectives) are provides the <provenance> tag, permitting versioning of metadata across systems, and are provides the <provenance> tag, permitting versioning of metadata across systems, and are provides the <provenance> tag, permitting versioning of metadata across systems, and are provides the <provenance> tag, permitting versioning of metadata across systems, and are provided in various contexts, such as the DART project (Dahlquist effective) are provided to the state of the st to a known identity is nes the sprovenances tag, permitting versioning of metadota across sprems. The been further refined in various contexts, such as the DART project (Dahlquist data in been further refined in various contexts, such as the metadota monitor's data in been further nodel refins on the accuracy of the metadota monitor's data in the model refins on the accuracy of the metadota monitor's data in the model refins on the accuracy of the metadota monitor's data in the model refins on the accuracy of the metadota monitor's data in the model refins on the accuracy of the metadota monitor's data in the model refins on the accuracy of the metadota metadota in the model refins on the accuracy of the metadota metadota in the model refins on the accuracy of the metadota metadota in the model refins on the accuracy of the metadota metadota in the model refins on the accuracy of the metadota in the metadota in the model refins on the accuracy of the metadota in the metadota in the model refins on the accuracy of the metadota in t dating a signature at the model relies on the accuracy of the metadata provider's data; in the model relies on the accuracy of the metadata provider's data; in identity has signed th that the same ident However, if a key the number of intermediate organisations relatively low, composed mostly compromised an upt the responsibilities originating from work of must.



the paper : multiple expressions, manifestations and copies



the presentation : expression(s) or new scholarlyWork?



capturing this in DC

- the DCMI Abstract Model (DCAM) says what the descriptions look like
- it provides the notion of 'description sets'
- i.e. groups of related `descriptions'
- where each 'description' is about an instance of one of the entities in the model
- relationships and attributes are captured as metadata properties in the application profile



from model to profile

- the application model defines the entities and relationships
- each entity and its relationships are described using an agreed set of attributes / properties
- the application profile describes these properties
 - contains recommendations, cataloguing/usage guidelines and examples
 - little is mandatory, prescriptive statements are limited
 - structured according to the entities in the model



application profile and vocabularies



the application profile

- DC Metadata Element Set properties (the usual simple DC suspects ...)
 - identifier, title, abstract, subject, creator, publisher, type, language, format
- DC Terms properties (qualified DC)
 - access rights, licence, date available, bibliographic citation, references, date modified
- new properties
 - grant number, affiliated institution, status, version, copyright holder
- properties from other metadata property sets
 - funder, supervisor, editor (MARC relators)
 - name, family name, given name, workplace homepage, mailbox, homepage (FOAF)
- clearer use of existing relationships
 - has version, is part of
- new relationship properties
 - has adaptation, has translation, is expressed as, is manifested as, is available as
- vocabularies

UKOLN

access rights, entity type, resource type and status

example properties

UKOLN

					Agent: name	
ScholarlyWork: title				type of agent date of birth		
subject abstract	Expression: title date available status version number language genre / type copyright holder bibliographic cita identifier				mailbox homepage identifier	
institution identifier		Manifestati		ion:		
		ation	moainea	Cop date acces licent ident	y: available ss rights ce ifier	

oai-pmh, dumb-down and community acceptance



OAI-PMH, dumb-down

- dumb-down
 - we still need to be able to create simple DC descriptions
 - we have chosen to dumb-down to separate simple DC descriptions of the ScholarlyWork and each Copy
 - simple DC about the ScholarlyWork corresponds to previous guidance
 - simple DC about each Copy is useful for getting to fulltext, e.g. by Google

XML schema

- produced by Pete Johnston, Eduserv Foundation
- specifies an XML format (Eprints-DC-XML) for representing a DC metadata *description set*
- based closely on a working draft of the DCMI Architecture Working Group for an XML format for representing DC metadata (DCXMLFULL)
- enables the creation, exposure and sharing of Eprints DC XML (epdcx)



community acceptance

- community acceptance plan outlines further work towards community take-up
 - deployment by developers
 - deployment by repositories, services
 - dissemination
 - DC community *may* take forward development of the profile
- more application profiles
 - JISC is funding work on profiles for images, timebased media and geographic data
 - this approach may prove a good foundation



thoughts on the approach ...

- this approach is guided by the functional requirements identified and the primary use case of richer, more functional, metadata
- it also makes it easier to rationalise 'traditional' and 'modern' citations
 - traditional citations tend to be made between eprint 'expressions'
 - hypertext links tend to be made between eprint 'copies' (or 'items' in FRBR terms)
- a complex underlying model may be manifest in relatively simple metadata and/or end-user interfaces
- existing eprint systems may well capture this level of detail currently – but use of simple DC stops them exposing it to others!
- it is the DCAM that allows us to do this with Dublin Core



thank you!

Julie Allinson j.allinson@ukoln.ac.uk www.ukoln.ac.uk/repositories/digirep/

