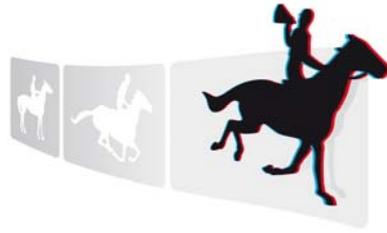


KEEPING AUDIOVISUAL CONTENT ALIVE



PRESTO
CENTRE

Metadata models, interoperability gaps and extensions to preservation metadata standards

Guus Schreiber

2010 Guus Schreiber

Published by PrestoCentre, <http://www.prestocentre.org>

The editors and authors of this work assert their moral rights to be identified as such.

This work may be copied, distributed and transmitted for all but commercial purposes. You may not alter, transform, or build upon this work. For any reuse or distribution, you must make clear to others the license terms of this work, including any embedded licensing metadata.

Enquiries concerning use of this work outside these terms should be sent to PrestoCentre, % Netherlands Institute for Sound and Vision, PO Box 1060, 1200 BB Hilversum, The Netherlands.

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivs 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/3.0/> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.

Table of contents

Scope.....	6
Executive summary.....	7
1 Expectations on (preservation) metadata standards with regard to a/v content.....	8
2 Models for metadata types in the digital library community and cultural heritage domain.....	9
2.1 Models for descriptive metadata.....	9
Dublin Core and the DCMI Metadata Terms.....	9
EAD.....	10
museumdat and LIDO.....	13
Europeana – the European Digital Library.....	14
2.2 Specific metadata models for a/v content.....	17
MPEG-7.....	17
P_Meta.....	18
PrestoSpace model.....	19
W3C Media Annotation WG.....	20
W3C Media Fragments WG.....	23
European Film Gateway metadata model.....	23
SMPTE.....	24
EBU Core.....	24
TV Anytime.....	27
VideoActive.....	28
VRA.....	29
PBCore.....	29
2.3 Models for preservation metadata.....	30
PREMIS.....	30
Intellectual Entities.....	30
Objects.....	30
Events.....	31
Rights.....	31
Agents.....	32
METS (wrapper format).....	32
Metadata model of the New Zealand National Library (based on PREMIS and METS).....	33
DNX.....	35
Long-term preservation Metadata for Electronic Resources (LMER).....	36

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

Shared format registries.....	36
2.4 Models for provenance metadata.....	37
ORE.....	37
PREMIS.....	37
METS.....	38
Open Provenance Model (OPM).....	38
Changeset.....	38
Provenance Vocabulary.....	38
Provenir.....	39
Content identifiers.....	39
2.5 Models for rights metadata.....	39
3 Shortcomings and gaps in metadata standards and models.....	40
3.1 Preservation metadata.....	40
3.2 Provenance metadata.....	41
3.3 Rights metadata.....	42
4 Conclusions and Recommendations.....	43
4.1 Integration of preservation and provenance metadata models.....	43
4.2 Relation to Europeana.....	43
Glossary.....	45
References.....	47

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

Scope

Audiovisual content collections are undergoing a transformation from archives of analogue materials to very large stores of digital data. PrestoPRIME researches and develops practical solutions for the long-term preservation of digital media objects, programmes and collections. This will be used for finding ways to increase access by integrating the media archives with European on-line digital libraries in a digital preservation framework.

The current report analyses existing models for different types of metadata used in the digital library community and elsewhere in the cultural heritage domain. Special attention was laid on descriptive and preservation metadata for identifying the shortcomings of such standards and models with regard to audiovisual content as dealt within PrestoPRIME. Also included in this report is summarising information on the use of metadata formats and standards. Resulting from this analysis an approach is proposed for which standards to use and where to extent them as needed. The approach will be the basis for implementation in the PrestoPRIME prototypes.

This report pays special attention to metadata interoperability with Europeana, in particular the new Europeana Data model, as Europeana has been identified as a major output channel for PrestoPRIME.

For analysing the needs on existing preservation standards (like PREMIS) to support audiovisual content a survey was created and several archives contacted.

Rights metadata are not within scope of this task and report. They are dealt with in WP4T4 and corresponding deliverable D4.0.5.

Executive summary

PrestoPRIME (<http://www.prestoprime.org/>) is an EU-funded project which focuses on developing practical solutions for the long-term preservation of digital media objects, programmes and collections, and find ways to increase access by integrating the media archives with European on-line digital libraries in a digital preservation framework. The project should deliver a range of tools and services, delivered through a networked Competence Centre.

The *purpose* of this document is to *layout* in detail the *metadata landscape* in which the project has to achieve its results. This landscape is extremely hybrid: there is an abundance of standards to consider, both from the broadcast and from the Web world. Some standards are proprietary in nature; other are open and free. PrestoPRIME has to deal with these different metadata cultures. We need to prevent at all costs that the project is developing “yet another standard”. This would only make the goal of achieving interoperability even harder to reach. Metadata interoperability between European broadcast archives as well as with other European collection holders is an essential requirement for PrestoPRIME.

This document therefore contains an extensive portfolio of metadata standards relevant for PrestoPRIME. We emphasize in particular the way in which metadata standards overlap, and where needed, provide extensive mapping tables for relations between standards.

With respect to sorts of metadata we take the following dimensions into account: (i) descriptive versus administrative metadata; (ii) AV-specific versus AV-neutral metadata, (iii) the nature of provenance metadata. One can see these distinctions in the structure of the document: each of the chapters is structured along one of these dimensions.

Special attention is given in this document to metadata interoperability with Europeana, as the project intends to be an important “gateway” for objects from broadcast archives to the Europeana portal. The document contains an in-depth discussion of this issue, taking also into account results from other projects that work with streaming video (e.g. EUScreen).

It should be noted that PrestoPRIME is not content with just this deliverable, but the members of the consortium are actually actively influencing the current activities in the relevant metadata standardization initiatives. For example, members of the consortium are actively involved in the W3C activities on Media Fragment and Metadata Annotation as well as in the Europeana Data Model.

1 Expectations on (preservation) metadata standards with regard to a/v content

Metadata are to be used for access and fruition of content (descriptive metadata) and for preservation purposes (technical metadata and identifying metadata).

The main difference for metadata in the audiovisual domain compared to those in the “classical” cultural heritage area exist due to the time dimension of a/v content. This starts with additional structuring and metadata (e.g. time codes). Further this affects other descriptive information as many elements can depend on temporal information ranging e.g. from annotations (on programmes, scenes, shots etc.) down to low level feature descriptors that may result from quality analysis operations on a/v content or maybe even from automated extraction from digital a/v material.

Existing metadata standards were analysed for how they can be used and what information may not be covered sufficiently within those standards. The survey included metadata representations for descriptive metadata (general models like Dublin Core, Encoded Archival Description and the Europeana Data Model as well as models especially intended for representation of a/v content like MPEG-7, P_Meta, PBCore and others).

Regarding preservation metadata several initiatives exist in Europe and in New Zealand. These and the standards PREMIS and METS were looked into. Provenance information is covered to some extent with the aforementioned standards. The representation of provenance information in further initiatives was also investigated.

Potential weaknesses in the analysed approaches resulted in a proposed integration of different information units within METS as a container.

2 Models for metadata types in the digital library community and cultural heritage domain

A number of metadata models and standards exist in the digital library and in the cultural heritage areas, covering different types of metadata.

2.1 Models for descriptive metadata

There are several general (i.e. not domain specific) metadata models applied across the digital library and cultural heritage area.

Dublin Core and the DCMI Metadata Terms

Rather general metadata models are in use covering several sub domains of cultural heritage. One of them is Dublin Core, in earlier times developed as the Dublin Core Metadata Element Set, Version 1.1¹ with 15 properties. These were core properties and therefore kept relatively general.

The Dublin Core model (maintained by the Dublin Core Metadata Initiative²) evolved over time. The full set of elements is an elaborate extension of the original 1.1 version and is called the DCMI Metadata Terms³.

A term is described by a number of attributes:

- **Name:** A token assigned to the term, unique within the term's DCMI namespace.
- **Label:** The human-readable label assigned to the term.
- **URI:** The Uniform Resource Identifier used to uniquely identify a term.
- **Definition:** A statement that represents the concept and essential nature of the term.
- **Type of Term:** The type of term as described in the DCMI Abstract Model [DCAM].

The DCMI Metadata Terms are used within the DCMI Abstract Model⁴. DCMI terms are specified using an RDF model, which defines the semantics of the DCMI elements. The DCMI terms form a hierarchy of descriptive properties, specified with the help of the `rdfs:subProperty` relation. For example, `dcterms:coverage` (an element of the original set of 15 DC elements) has two specializations, namely `dcterms:spatial` (for spatial coverage) and `dc:terms:temporal` (for temporal coverage).⁵

The question is: what about a/v content? How can this be represented with DCMI Metadata Terms? Two terms which could be used for structuring objects are ***hasPart*** and ***isPartOf***. But terms for further description of substructures (e.g. a time code) are not available. We come back to this issue in the section on the Europeana Data Model (EDM). EDM uses the DCMI terms but has a separate construct for representing part-of relations.

¹ <http://dublincore.org/documents/dces/> (verified 26 April 2010)

² <http://dublincore.org/> (verified 10 June 2010)

³ <http://dublincore.org/documents/dcmi-terms/> (verified 26 April 2010)

⁴ <http://dublincore.org/documents/abstract-model/> (verified 26 April 2010)

⁵ The syntax "`dcterms:coverage`" is a shorthand for the URI of the DCMI Terms namespace (`dcterms`) and the local identifier of the resource (i.e. `coverage`). This shorthand is called a "qname" and is now accepted practice in the Web world.

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

EAD

The Encoded Archival Description (EAD⁶) is an XML standard. It was developed within and for the archives' world and is also partially used in museums and libraries. Its main use is to exchange data of archival finding aids. The development started with a representation in SGML. This was later on changed to DTD and resulted (2002) in an XML schema.

According to the EAD Design principles (the way proposed to enhance EAD in the future⁷) EAD is a data structure but not a data content standard. Therefore it is not prescribed how an archive formulates the data appearing in the data elements of a document. This is left for external national or international data content standards. In addition to the structure the so called EAD Tag Library is provided. The tag library illustrates the type of data to be included in the particular elements.

Nothing specific was found about the description of a/v material with EAD documents. On the other hand there were activities to harmonise or map descriptions in EAD on the data model developed in Europeana. The results are shown in a technical report on Archival Digital Object Ingestion into Europeana (ESE-EAD harmonisation) Version 1.0, 07/08/2009⁸.

Regarding the description of a/v content there was an endeavour from the Archives of American Art at the Smithsonian Institution⁹. The resulting document "*Examples of EAD encoding for the description of audiovisual materials*"¹⁰ is not a recommendation on best practices but it gives an insight into today's practical approaches dealing with a/v content description in EAD schema as done in several archives. Three examples could be found therein with descriptions of the structure of a/v content. The other cases only describe the physical items used to store a/v content.

The following three examples show different aspects of the time dimension being part of a/v content descriptions. Different ways, but still not all of these descriptions provide time codes as would be necessary.

⁶ <http://www.loc.gov/ead/index.html> (verified 26 April 2010)

⁷ <http://www.loc.gov/ead/eaddesign.html> (verified 26 April 2010)

⁸ http://www.europeana-local.at/images/technical_report_archival_digital_object_ingestion_into_europeana_ese-ead_harmonisation_v1.0.pdf (verified 26 April 2010)

⁹ <http://www.aaa.si.edu/> (verified 26 April 2010)

¹⁰ <http://ead-for-av.googlegroups.com/web/EAD+XML+excerpts.doc?gda=hZowkEcAAACESXQ3HtQXbTEiPvdVDXbamHC03ot8ZihWYtbz62A3gdfaXpEJGv60pK0pqvug-PvdXO1PrwmeTtBIQmGaFzJKeV4duv6pDMGhhhZdjQINAw> (verified 26 April 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

Example excerpt from the Irish Virtual Research Library and Archive Project¹¹

```
<mods xsi:schemaLocation="http://www.loc.gov/mods/v3
http://www.loc.gov/standards/mods/v3/mods-3-2.xsd" version="3.2">
  <titleInfo><title>Recording</title></titleInfo>
  <physicalDescription><extent>00:19:39</extent>
    <reformattingQuality>preservation</reformattingQuality>
    <internetMediaType>audio/wav</internetMediaType>
    <digitalOrigin>reformatted digital</digitalOrigin>
  </physicalDescription>
  <note displayLabel="DVD number" type="admin">IVRLA_60005</note>
  <identifier type="local">IF_PM_6000062</identifier>
  <location><physicalLocation type="repository">UCD School of Irish,
Celtic Studies, Irish Folklore and
Linguistics</physicalLocation><physicalLocation
type="originalRef">UFP0049</physicalLocation>
</location>
  <part type="IVRLAObject">
    <detail><number>OB_6000057_IF</number></detail>
  </part>
  <part type="audioVisualContents">
    <detail><number>00:00:00</number><caption>Fishing with his
friend and his brother</caption></detail>
  </part>
  <part type="audioVisualContents">
    <detail><number>00:01:28</number><caption>Seine nets, fishing
on the banks</caption></detail>
  </part>
  <part
type="audioVisualContents"><detail><number>00:02:00</number><capti
on>Fishing boat and fishing methods</caption></detail>
</part>
```

¹¹ <http://ivrlaprod.ucd.ie/fedora/get/ivrla10-:4635/ivrla10-:objLayoutbDef/getLayout/> (verified 26 April 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

[Interview with Larry Fitzpatrick [Fisherman ?], Ringsend, Dublin. Collected by Séamus Sisk.] Use this identifier to cite or link to this digital object: http://hdl.handle.net/10151/OB_6000057_IF

<p>Collection: Urban Folklore Project (Dublin) [10 items] </p> <p>Extent: 1 item</p> <p>Resource Type: Sound Recording-Nonmusical</p> <p>Genre: Interview</p> <p>Note: National Folklore Collection MS number: 1980:</p> <p>Roles: - Fitzpatrick, Larry (Interviewee) - Ó Catháin, Séamas (Research team head) - Sisk, Séamus (Researcher)</p> <p>Language(s): English</p> <p>Publication/Creation Place: Ringsend, Dublin 4, Ireland</p> <p>Creation Date(s): 11 January 1980</p> <p>Institution: UCD School of Irish, Celtic Studies, Irish Folklore and Linguistics</p> <p>Repository: UCD Delargy Centre for Irish Folklore and the National Folklore Collection</p> <p>Reference Number: UFP0049</p>	<p>Subjects:</p> <p>Fishing boats Fishing--Dublin (Ireland) Gouldings Fertilizer Company Guano Interviews--Dublin (Ireland) Lakes--Ringsend (Ireland) Seining Sewage disposal in rivers, lakes, etc--Dublin (Ireland) Shipwrecks--Dublin (Ireland)</p>
--	---

1. **Recording** 00:19:39 © 2009 Prowayer Ltd

Time:	Topic:
00:00:00	Fishing with his friend and his brother
00:01:28	Seine nets, fishing on the banks
00:02:00	Fishing boat and fishing methods
00:06:55	Fishing in Dublin Bay off Monkstown "the Image Gardens", Sutton and Dollymount
00:07:57	Wrecks in Dublin Bay
00:09:24	Guano
00:10:46	Gouldings fertilizer company on the East Wall
00:11:39	Unloading guano
00:14:35	Lake on Ringsend strand, the "Cross Lake"
00:15:24	The Costellos, job was to open and shut the sluice gates for discharging the sewage into the river
00:16:33	"The Shamrock" sewage ship

Figure 1: Example screen from the "Irish Virtual Research Library and Archive Project"

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

Example excerpt¹² from the *Archivo Luce in Italy*¹³ (Italian language only)

```

<physdesc label="video">
  <archref>
    <container type="scatola">
      <emph>Ex 52114/vt; Provenienza: lavorazione</emph>
      1
    </container>
    <unitid type="inventario">A/BETA/945</unitid>
  </archref>
  <dimensions unit="hh:mm:ss">00:18:00</dimensions>
  <extent label="completo" type="cassetta">1</extent>
  <genreform type="master">
    <emph>4/3</emph>
    BETA SP
  </genreform>
  <physfacet type="colore">b/n</physfacet>
  <physfacet type="sonoro">
    <emph>mono</emph>
    sonoro
  </physfacet>
</physdesc>
<physdesc label="video">
  ...
</physdesc>

```

Example excerpt from the “Register of the N. N. Poppe sound recording”¹⁴ at the Hoover Institution Archives

```

<dsc type="in-depth" id="dsc-1">
  <head>Collection Contents</head>
  <c01>
    <did>
      <container type="disc" label="Disc No./Sides: "> 1 :
      1,2</container>
      <unittitle>Disc recording of N. N. Poppe on the Academy of
      Sciences parts 1 and 2</unittitle>
      <physdesc>26:18 minutes </physdesc>
    </did>
    <scopecontent><p>Poppe speaks on the import and membership of the
    academy. He describes life in the academy, usually focusing on

```

¹² http://www.regesta.com/xdams/ontologie/xml/Archivio_Audiovisivo.xml (verified 26 April 2010)

¹³ <http://www.archivioluce.com/archivio/> (verified 26 April 2010)

¹⁴ http://www.oac.cdlib.org/view?docId=kt9779s0x7%3bdeveloper=local%3bquery=poppe%3bstyle=oac4%3bdoc.view=entire_text#hitNum5 (verified 26 April 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

cost, in the areas of various canteens, foreign books/literature, commuting, medical care, and the living space of his family; The disc begins with a short explanation of the cutting/recording process.</p></scopecontent>

</c01>

museumdat and LIDO

A further format developed within the cultural heritage area is museumdat¹⁵. The format was developed for publication of museum core data and research in museum portals. The (XML) scheme provides a number of attributes for core museum data and for some of them references to the authority files (e.g. Art and Architecture Thesaurus, German Name Authority File, Thesaurus of Geographic Names) proposed for use. The format was elaborated in a German working group called "Fachgruppe Dokumentation des Deutschen Museumsbundes".

A follow up exercise on a more international level is called Lightweight Information Describing Objects (LIDO). The "Data Harvesting and Interchange Working Group"¹⁶ supports the development of LIDO resulting from museumdat, CDWA lite¹⁷ (US; core records for works of art and material culture assisting union cataloguing), CIDOC CRM¹⁸ (overall conceptual model for use in cultural heritage documentation) and SPECTRUM¹⁹ (UK; collection management and exchange between collections).

Areas covered in the LIDO scheme include: descriptive and administrative metadata, identification and classification, actors, location, reproduction information, events and rights information.

Europeana – the European Digital Library

Initial version: ESE - Europeana Semantic Elements

Europeana is an initiative to set up the European cultural heritage portal. Europeana is being realised in a series of Europeana projects. The first version of Europeana used an interoperable metadata model called the Europeana Semantic Elements (ESE specifications Version 3.2²⁰, 07/08/2009). The Europeana Semantic Elements are used for maintaining the Europeana portal (functional specification²¹) and allowing access to a wide range of heterogeneous collections in the cultural heritage domain. Data from those collections are ingested into the central database of Europeana and made available to the public through a web user interface.

The Europeana Semantic Elements (ESE) are defined using the Dublin Core (DC) metadata elements, a subset of the DC terms and a set of elements which were created to meet Europeana's needs. The table below shows the element set:

¹⁵ <http://www.museumdat.org/> (verified 10 June 2010)

¹⁶ [http://cidoc.mediahost.org/WG_Data_Harvesting\(en\)\(E1\).xml](http://cidoc.mediahost.org/WG_Data_Harvesting(en)(E1).xml) (verified 10 June 2010)

¹⁷ http://getty.art.museum/research/conducting_research/standards/cdwa/cdwalite.html (verified 10 June 2010)

¹⁸ <http://www.cidoc-crm.org/> (verified 10 June 2010)

¹⁹ <http://www.collectionstrust.org.uk/spectrum>

²⁰ http://www.europeana-local.at/images/europeana_semantic_elements_specifications_v3.2.pdf (verified 26 April 2010)

²¹ <http://abm.ylm.se/europeanalocal/pdf/EuropeanaOutline08.pdf> (verified 26 April 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

Source	Element	Element Refinement(s)
DC	title	alternative
DC	creator	
DC	subject	
DC	description	tableOfContents
DC	publisher	
DC	contributor	
DC	date	created; issued
DC	type	
DC	format	extent; medium
DC	identifier	
DC	source	
DC	language	
DC & Europeana	relation	isVersionOf; hasVersion; isReplacedBy; replaces; isRequiredBy; requires; isPartOf; hasPart; isReferencedBy; references; isFormatOf; hasFormat; conformsTo; isShownBy; isShownAt
DC	coverage	spatial; temporal
DC	rights	
DC	terms provenance	
Europeana	userTag	
Europeana	unstored object	
Europeana	language	
Europeana	provider	
Europeana	type	
Europeana	uri	
Europeana	year	
Europeana	hasObject	
Europeana	country	

Table 1: Europeana Semantic Elements

New version: EDM- Europeana Data Model

Within the Europeana v1.0 project a new metadata model is under development. This Europeana Data Model (EDM) is scheduled to be used in future versions of Europeana. The outlook for EDM is that it should be finalised in the summer of 2010. For the mapping of metadata catalogues to EDM an ingestion platform is under development – the Athena ingest tool²².

The main rationale for developing EDM as a replacement for ESE is that ESE represents the lowest common denominator for Europeana object metadata. ESE reduces the original collection metadata to a subset of Dublin Core (see previous section). This forces interoperability, but at a high price. The major drawback is that the original metadata is lost. Institutions participating in Europeana have expressed their concerns about this simplification. The goal of EDM is therefore to preserve the original metadata while still

²² http://athena.image.ece.ntua.gr/athena/Login_input.action (verified 26 April 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

allowing for interoperability. This dual goal is achieved in EDM through a Semantic Web representation of EDM with the help of RDF²³.

The following requirements were formulated for EDM:

1. Distinction between “provided object” (painting, book, program) and digital representation.
2. Distinction between object and metadata record describing an object.
3. Allow for multiple records for same object, containing potentially contradictory statements about an object.
4. Support for objects that are composed of other objects.
5. Standard metadata format that can be specialized.
6. Standard vocabulary format that can be specialized.
7. EDM should be based on existing standards.

The requirements fit well with PrestoPRIME; in particular requirement 4 is important for the project. We therefore treat EDM in some detail in this document. At the time of writing the best information source is the EDM Primer²⁴.

EDM consists of three standards (see the last requirement) which together fill the requirements:

- Dublin Core is used for metadata representation (Requirement 5).
- SKOS for vocabulary representation (Requirement 6).
- OAI ORE is used for organization of metadata about an object (Requirements 1-4).

EDM uses the latest version of DCMI Metadata Terms²⁵. As remarked before the DCMI terms are a specialisation of the 15 original DC elements. The set is specified with an RDF model, which means that the DCMI terms can be specialised themselves. This makes it possible to define in the ingestion process a collection-specific metadata element as specialisation of a DCMI term, this fulfilling the overall EDM goal of not losing the original data.

SKOS²⁶ is a W3C standard for publication of vocabularies on the Web. It is used by large institutions. For example, Library of Congress has published all their Subject Headings in SKOS format. Similar to Dublin Core it is defined with an RDF model, allowing for vocabulary-specific specialisations of SKOS. The EuropeanaConnect project is setting up a vocabulary server with a large number of cultural-heritage and general-purpose vocabularies. Sound & Vision has their in-house GTAA vocabulary available in SKOS format²⁷.

OAI ORE²⁸ (Open Archives Initiative Object Reuse & Exchange) is a standard of the Open Access Initiative. Specification: it is also specified with an RDF model, allowing collection-

²³ <http://www.w3.org/TR/rdf-primer/> (verified 2 June 2010)

²⁴ http://www.few.vu.nl/~aisaac/edm/EDM_Primer_100401.pdf (verified 1 June 2010)

²⁵ <http://dublincore.org/documents/dcmi-terms/> (verified 2 June 2010)

²⁶ <http://www.w3.org/TR/skos-primer/> (verified 2 June 2010)

²⁷ <http://thesauri.cs.vu.nl/eswc06/> (verified 24 June 2010)

²⁸ <http://www.openarchives.org/ore/1.0/toc.html> (verified 2 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

specific specialisation of the ORE constructs. EDM deploys four ORE constructs (represented as RDF classes)

1. **Object:** the book/painting/program/video being described. Objects can be “digitally born” and should be distinguished from their digital (re-)representations (see item 3 below), e.g. in different formats or resolution.
2. **Aggregation:** organises object information from a particular provider (museum, archive and library). Provenance metadata are attached to the aggregation. For example, a `dcterms:creator` property on an aggregation might say that this object is created by INA. For one object there may be multiple aggregations, meaning that the object is part of multiple collections (e.g. BBC and RAI having the same program in their archive).
3. **Digital representation:** some digital form of the object with a Web address. There may be multiple digital representations per object (e.g. different AV formats).
4. **Proxy:** the metadata record for the object. This is the descriptive metadata of the object, e.g. the director of the movie, the subject of the movie, etc.

In addition, ORE specifies relations between these constructs.

The EDM Primer²⁹ shows extensive examples of how Dublin Core, SKOS and OAI ORE can be used together to represent Europeana objects. The fact that it caters for objects consisting of smaller objects is essential for PrestoPRIME (e.g. to represent a fragment of a program or an episode of a series).

EDM contains some additional definitions such as predefined classes for person, place, time and event, but these are not essential for the discussion here.

2.2 Specific metadata models for a/v content

The “more general” standards in the cultural heritage domain are not very widely used for a/v content and obviously there are some good reasons for that. In the a/v domain a number of further specific standards are available and in use for several purposes. In the following some of them are discussed.

MPEG-7

The ISO/IEC standard *Multimedia Content Description Interface* [MPEG7] has been defined as a format for the description of multimedia content in a wide range of applications. MPEG-7 defines a set of description tools, called description schemes (DS) and descriptors (D). Descriptors represent single properties of the content description, while description schemes are containers for descriptors and other description schemes. The definition of description schemes and descriptors uses the *Description Definition Language* (DDL), which is an extension of XML Schema. MPEG-7 descriptions can be either represented as XML (textual format, TeM) or in a binary format (binary format, BiM).

An important part of MPEG-7 are the *Multimedia Description Schemes* (MDS), which provide support for the description of media information, creation and production information, content structure, usage of content, semantics, navigation and access, content organisation and user interaction. The structuring tools are very flexible and allow the description of content on different levels of granularity. In addition, the *Audio* and *Visual* parts define low- and mid-level descriptors for these modalities.

²⁹ http://www.few.vu.nl/~aisaac/edm/EDM_Primer_100401.pdf (verified 10 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

MPEG-7 allows the description of audiovisual content in a very detailed way. This is due to its comprehensive and flexible representation. However these very enabling possibilities also create some risks when describing content and especially when exchanging such descriptions based on that standard: a lot of semantics can be included and this also means adding interpretation to descriptions; ambiguities can be added if not carefully designed and descriptions may become very complex.

The concept of profiles has been introduced to define subsets of the comprehensive standard which target certain application areas. Three profiles have been standardised: the *Simple Metadata Profile* (SMP), which describes single instances or collections of multimedia content, the *User Description Profile* (UDP), containing tools for describing personal preferences and usage patterns of users of multimedia content in order to enable automatic discovery, selection, personalisation and recommendation of multimedia content, and the *Core Description Profile* [CDP], which consists of tools for describing general multimedia content such as images, videos, audio and collections thereof.

One of the profiles that has been proposed with the goal of solving semantic interoperability issues is the Detailed Audiovisual Profile³⁰ (DAVP) [BS06]. The given profile allows detailed description of single multimedia content entities. Included are functionalities allowing a comprehensive structural description of the content, possibilities for textual and semantic annotations as well as adding audio and visual feature descriptions. Application areas covered with DAVP are applications that deal with the analysis, description, retrieval, summarization and exchange of audiovisual content.

The NHK Metadata Production Framework [MPF] data model is an industrial application of the Core Description Profile. The authors address the complexity and ambiguity problems of MPEG-7 proposing a metadata model that further restricts CDP by excluding some elements and reducing the cardinality of others. The new version also allows the use of the visual and audio descriptors defined in parts 3 and 4. The definition of the data model defines a number of semantic constraints for the structure of the description as well as several syntactic and semantic constraints on different elements of the description (called "operational rules").

Currently the EBU EC-M/SCAIE group³¹ is working on the definition of a profile to be used for describing results of automated information extraction tools in broadcast production processes. Both MPF and DAVP are used as inputs for the development of the new profile. The profile will be submitted as a proposal to the 93th MPEG meeting in July 2010.

P_Meta

P_META (EBU tech 3295³², available as version 2.1 since July 2009) was originally designed to support business to business content exchange (i.e. offering a standard vocabulary for information relating to programme information in the professional broadcasting industry).

- a universal standard for metadata exchanges between professional media organizations;
- a definition of common meaning to the data fields and values that most broadcasters use in order to

³⁰ <http://mpeg-7.joanneum.at/> (verified 26 April 2010)

³¹ <http://tech.ebu.ch/groups/pscaie> (checked 24 June 2010)

³² <http://tech.ebu.ch/docs/tech/tech3295v2.pdf> (checked 24 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

- enable exchanges;
- designed to be flexible and suitable for use in a wide range of broadcasting activities;
- both language and system independent;
- a joint development by EBU members on a not-for-profit basis;
- a scheme that makes use of other standards where possible, e.g. ISO country codes.

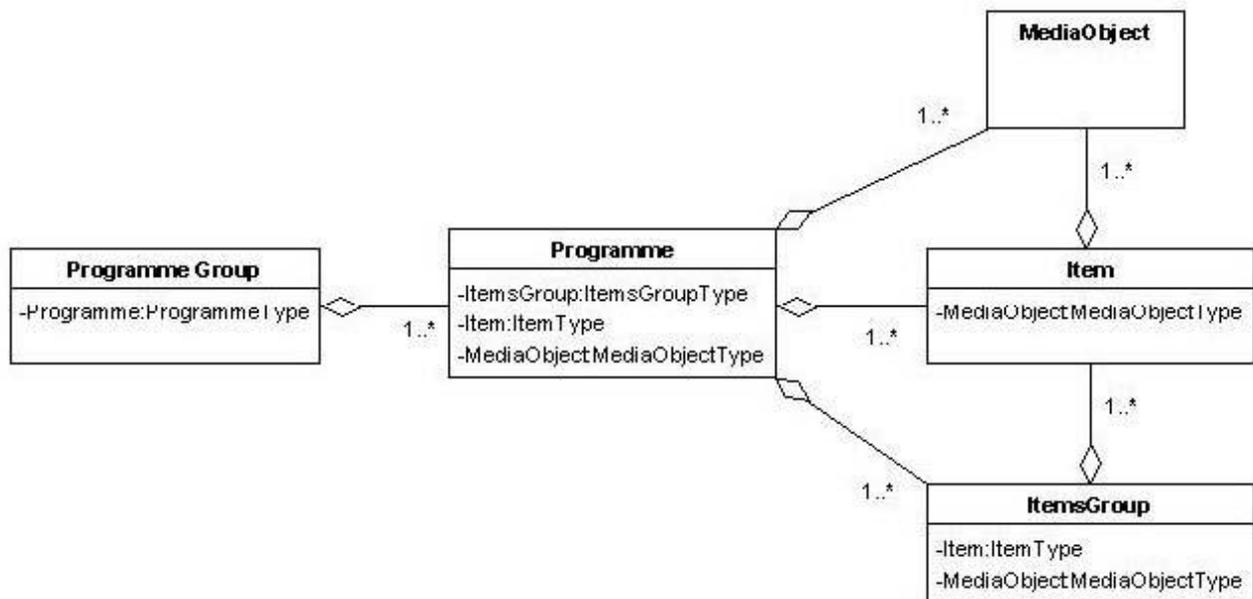


Figure 2: rough content structure around which P_META has been built

PrestoSpace model

The PrestoSpace project³³ document formats were defined³⁴ in order to build a framework for the documentation of audiovisual works (supporting processes including Preservation, Restoration, and Archive interfacing). Included in such composite XML documents – the Editorial Object Documents (EDOB) – are individual components expressed using the syntactic tools defined in external standards (e.g. P_META EBU Tech 3295, MPEG-7) that are most appropriate for the particular tasks.

³³ <http://prestospace.org> (verified: 10 May 2010)

³⁴ <http://www.crit.rai.it/attivita/PrestoSpaceFormats/PrestoSpaceFormats.html> (verified: 10 May 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

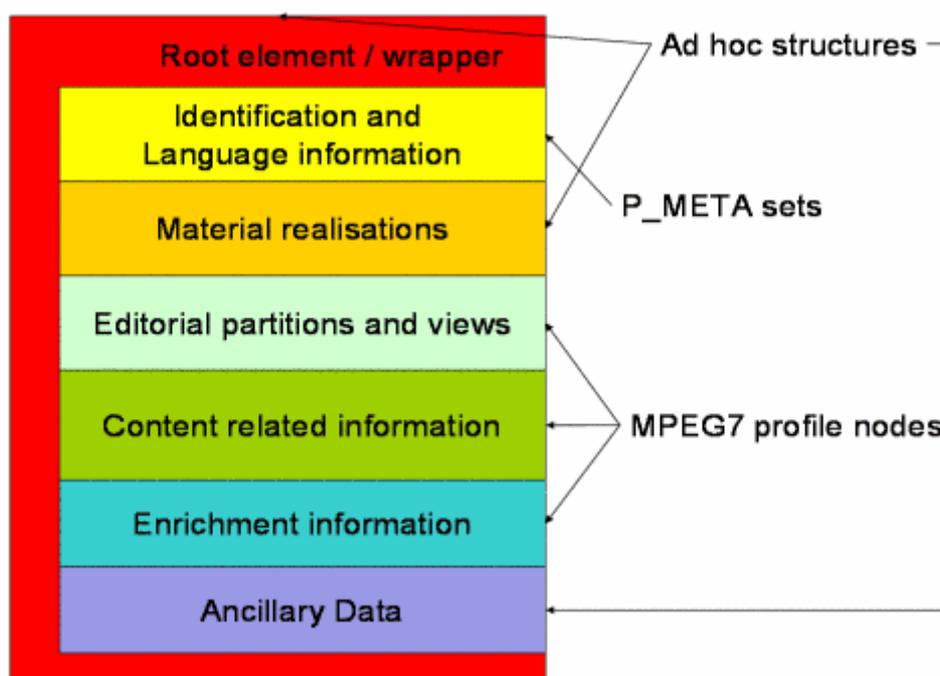


Figure 3: schematisation of the metadata format structure of PrestoSpace

W3C Media Annotation WG

The Media Annotations Working Group (MAWG) of the World Wide Web Consortium³⁵ (W3C) develops the so called Ontology for Media Resource (version 1.0 is an editors' copy without official standing yet)³⁶. It is supposed to become a core vocabulary to describe media resources on the Web and covers basic metadata items. The further defined semantics shall preserve mappings between existing formats. With the amount of interoperability introduced a certain loss of information is possible when mapping between formats.

A number of formats were selected by the working group for which mappings will be made available. They are shown in the following table from the website:

Identifier	Format	Example	Reference
cl11	CableLabs 1.1	cl11:Writer_Display	Cablelabs 1.1
cl20	CableLabs 2.0	cl20:Producer	Cablelabs 2.0
dig35	DIG35	dig35:ipr_name/ipr_person@description='Image Creator'	DIG35
dc	Dublin Core	dc:creator	Dublin Core
ebucore	EBUCore	ebuc:creator	EBUCore
pmeta	EBU P-Meta	pmeta:Contribution	EBU P-META
exif	EXIF 2.2	exif:Artist	EXIF
frbr	FRBR	frbr:Person	FRBR
id3	ID3	id3:TCOM	ID3
iptc	IPTC	iptc:Creator	IPTC
it	iTunes	it:©ART	iTunes

³⁵ <http://www.w3.org/2008/WebVideo/Annotations/> (verified: 17 December 2009)

³⁶ <http://www.w3.org/TR/2010/WD-mediaont-10-20100608/> (verified: 11 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

Identifier	Format	Example	Reference
lom21	LOM 2.1	lom21:LifeCycle/Contribute/Entity	LOM
ma	Core properties of MA WG	ma:creator	Property definition
media	Media RDF	media:Recording	Media RDF
mrss	Media RSS	mrss:credit@role='author'	Media RSS
mets	METS	mets:agency	METS
mpeg7	MPEG-7	mpeg7:CreationInformation/Creation/Creator/Agent	MPEG-7
nmix	NISO MIX	nmix:ImageCreation/ImageProducer	MIX
qt	Quicktime	qt:@dir	QuickTime
media	SearchMonkey Media	media:type	MediaMonkey
dms	DMS-1	dms:Participant/Person	DMS-1
tva	TV-Anytime	tva:CredistsList/CredistItem	TV-Anytime
txf	TXFeed	txf:author	TXFeed
vra40	VRA Core 4.0	vra40:agent	VRA
xmp	XMP	xmpDM:composer	XMP
yt	YouTube Data API Protocol	yt:author	YouTube Data API Protocol

Table 2: formats for mapping on W3C Ontology for Media Resource

The core properties of the media ontology³⁷ are listed below.

MAWG	Description	Media	Example
ma:identifier	A tuple identifying a resource, which can be either an abstract concept (e.g., Hamlet) or a specific object, using a URI. The type can be used to optionally define the category of the identifier.	all	http://www.w3.org/2008/WebVideo/Annotations/wiki/Image:MAWG-Stockholm-20090626.JPG
ma:title	A tuple providing the title or name given to the resource. The type can be used to optionally define the category of the title.	all	"MAWG-Stockholm-20090626"
ma:language	The language used in the resource. Recommended best practice is to use a controlled vocabulary such as [BCP 47].	all	en-us
ma:locator	The address at which the resource can be accessed (e.g. a URL, or a DVB URI).	all	http://www.w3.org/2008/WebVideo/Annotations/wiki/images/9/93/MAWG-Stockholm-20090626.JPG
ma:contributor	A tuple identifying the agent (with either a URI, if it exists, or plain text) and the nature of the contribution, e.g. actor, cameraman, director, singer, author, artist.	all	{imdb:nm0000318, director}
ma:creator	The author of the resource and the role. The author identifier can be defined as either an URI (which is best practice) or as plain text. The role is defined as plain text.	all	{dbpedia:Shakespeare, playwright}

³⁷ <http://www.w3.org/2008/WebVideo/Annotations/drafts/ontology10/WD/summary.html>
(verified: 15 December 2009)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

MAWG	Description	Media	Example
ma:createDate	The date defines the date and time that the resource was created. The type defines the particular category of creation date (e.g., release date, date recorded, date edited).	all	2009-06-26T15:30:00
ma:location	A location name and/or data where the resource has been shot/recorded.	all	" Stockholm, Kistavägen 25, KISTA, KMUM Building"
ma:description	Free-form text describing the content of the resource.	all	"Group picture of the W3C Media Annotations WG at the face-to-face meeting in Stockholm."
ma:keyword	A concept, descriptive phrase or keyword that specifies the topic of the resource. A recommended best practice is to take this keyword from an ontology or a controlled vocabulary.	all	"W3C Media Annotations WG"
ma:genre	The category of the content of the resource. Recommended best practice is to use an ontology or a controlled vocabulary such as the EBU vocabulary.	all	"work"
ma:rating	A tuple defining the rating value, the rating person or organization (as a URI or a string), and the voting range (min. value, max. value).	all	{[http://www.individuals.com/ChrisPope, 10.0, 0, 10.0, "quality"] (Rating person: http://www.individuals.com/ChrisPope, Rating value: 10.0, Rating min: 0, Rating max: 10.0, Rating context: "quality")}
ma:relation	A tuple identifying a resource to which the current resource is related and optionally, the nature of the relationship. An example is a listing of content that has a relationship (possibly a named) to another content.	all	{http://www.w3.org/2008/WebVideo/Annotations/wiki/Image:MAWG-Stockholm-20090626_thumb.JPG, "thumbnail"}
ma:collection	The URI (best practice) or the name of the collection from which the resource originates or to which it belongs.	all	"My Work Pictures"
ma:copyright	The copyright statement associated with the resource and optionally, the identifier of the copyright holder. Other issues related to Digital Rights Management are out of scope for this specification.	all	{"All images in the collection are copyrighted by Wonsuk Lee", http://www.individuals.com/WonsukLee}
ma:policy	A description of the security policy applying to the media resource, or a reference to the security policy (e.g., Creative Commons). The type attribute can be used to provide more information as to the nature of the security policy (e.g., permissions, access control, ownership).	all	{"Attribution 2.5 ", http://www.organizations.com/CreativeCommons}
ma:publisher	The publisher of a resource.	all	http://www.individuals.com/WonsukLee
ma:target Audience	A tuple identifying the issuer of the classification (parental guidance issuing agency, targeted geographical region) and the value given in this classification.	all	[http://www.fosi.org/icra,"no nudity"]
ma:fragment	A tuple containing a fragment identifier and its role. A fragment is a portion of the resource, as defined by the [MediaFragment] Working Group.	all	{"Person", http://www.example.com/movie.mov#xywh=320,320,40,100}
ma:named Fragment	A tuple containing a named fragment identifier and its label.	all	{"Joakim Söderberg", http://www.w3.org/2008/WebVideo/Annotations/wiki/Image:MAWG-Stockholm-

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

MAWG	Description	Media	Example
			20090626.JPG#xywh=1600,550,80,150}
ma:frameSize	The frame size of the resource, if applicable. For example: w:720, h:480. It is optional to specify the units; the default value is pixels.	I, V	{3.072, 2.304}
ma:compression	The compression type used. For container files (e.g., QuickTime, AVI), the compression is not defined by the format, as a container file can have several tracks with different encodings. In such a case, several ma:compression instances will exist. Thus, querying the ma:compression property of the track media fragments will return different values for each track fragment. Note: it is possible to use an extended MIME type as the value for this property, see [RFC 4281].	I, V, A, T	"jpeg"
ma:duration	The actual duration of the resource. The unit is defined to be seconds.	V, A	2.134
ma:format	The MIME type of the resource (e.g., wrapper, bucket media types).	I, V, A, T	"image/jpeg"
ma:samplingrate	The audio sampling rate. The unit is defined to be samples/second.	A	44100
ma:framerate	The video frame rate. The unit is defined to be frames/second.	V	25
ma:averageBitrate	The average bit rate. The unit is defined to be kbps.	A, V	4000
ma:numTracks	The number of tracks of a resource, optionally followed by the type of track (e.g., video, audio, subtitle).	A, V, T	{2,"audio"}

Key to media types: I ... Image; V ... Video; A ... Audio; T ... Text (e.g. closed caption)

Table 3: Core properties of the W3C Ontology for Media Resource

There was a last call on the current Working Draft in mid June 2010. After processing and replying to all requests following on that call the document may become a Candidate Recommendation in Fall 2010 leading to a Recommendation foreseen for early 2011. Besides the ontology the "API for Media Resource 1.0"³⁸ was specified to access elements referring to the ontology.

W3C Media Fragments WG

A major aspect when dealing with audiovisual content and its description is the fragmentation of content and its time dependency. These aspects have to be also covered within the persistency of addresses for fragments. The work done by the Media Fragments Working Group is supposed to be the answer to such questions:

*The mission of the Media Fragments Working Group, part of the Video in the Web Activity, is to address temporal and spatial media fragments in the Web using Uniform Resource Identifiers (URI).*³⁹

Several aspects are covered within the current W3C Working Draft 13 April 2010 of the *Media Fragments URI 1.0*⁴⁰. The specification is expected to go to last call working draft in summer 2010.

³⁸ <http://www.w3.org/TR/2010/WD-mediaont-api-1.0-20100608/> (verified: 11 Jun 2010)

³⁹ <http://www.w3.org/2008/WebVideo/Fragments/> (verified: 10 May 2010)

⁴⁰ <http://www.w3.org/2008/WebVideo/Fragments/WD-media-fragments-spec/> (verified: 10 May 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

European Film Gateway metadata model

The Project European Film Gateway⁴¹ developed an interoperability schema for the content providers represented in the project. The report [DEBOLE2009] describes the schema to be used for archival resources and filmographic descriptions. A further aim was to get a mapping from this schema to the Europeana Semantic Elements (ESE).

For object description three levels are in use: creation (referring to the concept Cinematographic Work), Manifestation (a/v manifestation for films and non-a/v manifestation for pictures, photos, books and so forth include all properties of the digital representation which may change over lifetime of the creation object but which do not affect the identity of the object) and Item (logical wrapper for the digital file). The model foresees eight major entities (a/v and non a/v creation, a/v and non a/v manifestation, collection, item, event and agent) which exist on the three above mentioned levels.

Temporal segmentation of content is mentioned but not specified. Instead it is stated that parts of MPEG-7 could be applied.

SMPTE

The Material Exchange Format [MXF] is a standard issued by Society of Motion Picture and Television Engineers (SMPTE), defining the specification of a file format for the wrapping and transport of essence and metadata in a single container. The Material Exchange Format is an open binary file format targeted at the interchange of captured, ingested, finished or “almost finished” audio-visual material with associated data and metadata. Support for technical metadata is built directly into MXF specification. In order to provide enough flexibility to deal with different kinds of descriptive metadata, a plugin mechanism for descriptive metadata is defined. These descriptive metadata schemes (DMS) can be integrated into MXF files. So far SMPTE has standardised the Descriptive Metadata Scheme 1 (DMS-1) and the EBU has defined a DMS for P_Meta.

The SMPTE Descriptive Metadata Scheme 1 (DMS-1, formerly known as Geneva Scheme [DMS-1]) uses metadata sets defined in the SMPTE Metadata Dictionary (see below). Metadata sets are organised in descriptive metadata (DM) frameworks. DMS-1 defines three DM frameworks, which correspond to different granularities of description: production (entire media item), clip (continuous AV essence part) and scene (narratively or dramatically coherent unit). When DMS-1 descriptions are embedded into MXF files they are represented in KLV (Key-Length-Value) format, but there exists also a serialised format based on XML Schema.

The SMPTE Metadata Dictionary [RP210] is not a metadata format on its own, but a large thematically structured list of narrowly defined metadata elements, defined by a key, the size of the value and its semantics. It is used for all metadata embedded in MXF files, but the elements defined in the dictionary are also used outside MXF. A common use is for metadata in the headers of DPX files. The dictionary has a comprehensive list of technical and process-related metadata elements.

EBU Core

The “EBU Core” set of metadata is said to be the minimum information needed to describe radio and television content (from the EBU Core specification [EBU2009]).

⁴¹ <http://www.europeanfilmgateway.eu/index.php> (verified: 11 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

The EUscreen⁴² project will provide access to distributed audiovisual heritage with mechanisms built on EBU Core and open web standards.

There already exists a mapping⁴³ between EBU Core and MAWG as shown in the table below:

MAWG	EBU Core	How to do the mapping	Datatype	XPath
Note: each is to be preceded by the namespace 'ebucore:'				
ma:contributor	contributor	Either a person or an organisation	string	EBUCoreMain/coreMetadata/contributor/contactDetails/name/name EBUCoreMain/coreMetadata/contributor/organisationName
ma:creator	creator	Either a person or an organisation	string	EBUCoreMain/coreMetadata/creator/contactDetails/name/name EBUCoreMain/coreMetadata/creator/organisationName
ma:description	dc:description	Free text	string	EBUCoreMain/coreMetadata/description/dc:description
ma:format	dc:format	Free text or a series or more specific attributes provided in the XPath column, which would qualify to map into ma:format	string	EBUCoreMain/coreMetadata/format/dc:format EBUCoreMain/coreMetadata/format/medium/@typeLabel or /@typeLink EBUCoreMain/coreMetadata/format/mimeType/@typeLabel or /@typeLink EBUCoreMain/coreMetadata/format/fileFormat/@typeLabel or /@typeLink etc.
ma:identifier	dc:identifier	DC compliance requires a string but it is recommended to use URIs or IRIs instead	string	EBUCoreMain/coreMetadata/identifier/dc:identifier
ma:language	dc:language languageCode	A free text term and/or a reference to a web resource such as a classification scheme term	string anyURI	EBUCoreMain/coreMetadata/language/dc:language EBUCoreMain/coreMetadata/language/languageCode
ma:publisher	dc:publisher	Either a person or an organisation	string	EBUCoreMain/coreMetadata/publisher/contactDetails/name/name EBUCoreMain/coreMetadata/publisher/organisationName
ma:relation	dc:relation dc:identifier relationLink	Free text or an identifier or a link to a related resource. Specialised relations are provided in EBU which would qualify for mapping (see XPaths)	string string anyURI	EBUCoreMain/coreMetadata/relation/dc:relation EBUCoreMain/coreMetadata/relation/relationIdentifier/dc:identifier EBUCoreMain/coreMetadata/relation/relationLink same apply to EBUCoreMain/coreMetadata/isVersionOf EBUCoreMain/coreMetadata/hasVersion EBUCoreMain/coreMetadata/isReplacedBy EBUCoreMain/coreMetadata/replaces EBUCoreMain/coreMetadata/isRequiredBy EBUCoreMain/coreMetadata/requires EBUCoreMain/coreMetadata/ispartOf EBUCoreMain/coreMetadata/hasPart EBUCoreMain/coreMetadata/isreferencedBy EBUCoreMain/coreMetadata/references EBUCoreMain/coreMetadata/isFormatOf EBUCoreMain/coreMetadata/hasFormat

⁴² <http://www.euscreen.eu/> (checked: 24 June 2010)

⁴³ <http://www.w3.org/2008/WebVideo/Annotations/drafts/ontology10/WD/EBUCore.html> (checked: 17 December 2009)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

MAWG	EBU Core	How to do the mapping	Datatype	XPath
Note: each is to be preceded by the namespace 'ebucore:'				
ma:keyword	dc:subject subjectCode	A free text term and/or a reference to a web resource such as a classification scheme term	string anyURI	EBUCoreMain/coreMetadata/subject/dc:subject EBUCoreMain/coreMetadata/subject/subjectCode
ma:title	dc:title	title/dc:title and/or alternativeTitle/dc:title	string	EBUCoreMain/coreMetadata/title/dc:title EBUCoreMain/coreMetadata/alternativeTitle/dc:title
ma:genre	dc:type genre/@type Label genre/@type Link	Either a free text term in dc:type (not necessarily only genre) or genre/@typeLabel or a reference to a web resource such as a classification scheme term genre/@typeLink	string string anyURI	EBUCoreMain/coreMetadata/type/dc:type EBUCoreMain/coreMetadata/type/genre/@typeLabel EBUCoreMain/coreMetadata/type/genre/@typeLink
ma:createDate	created		date	EBUCoreMain/coreMetadata/date/created
ma:rating				EBUCoreMain/coreMetadata/
ma:collection	Title	The 'type' of content being described should be "collection" in type/objectType/@typeLabel (free text) or type/objectType/@typeLink (anyURI to refer to e.g. a classification scheme term)	string	EBUCoreMain/coreMetadata/title/dc:title + EBUCoreMain/coreMetadata/type/objectType/@typeLabel (collection, string) EBUCoreMain/coreMetadata/type/objectType/@typeLink (collection, anyURI)
ma:duration	duration	It is important to look at the format used for expressing the duration in duration/@formatLabel or duration/@formatLink	string	EBUCoreMain/coreMetadata/format/duration
ma:copyright	dc:rights rightsLink exploitationIssues	Free text or a link to a web page with rights declaration or more specifically exploitation issues	string anyURI string	EBUCoreMain/coreMetadata/rights/dc:rights EBUCoreMain/coreMetadata/rights/rightsLink EBUCoreMain/coreMetadata/rights/exploitationIssues
ma:license	dc:rights rightsLink		string anyURI	EBUCoreMain/coreMetadata/rights/dc:rights EBUCoreMain/coreMetadata/rights/rightsLink
ma:location	dc:coverage name code posx + posy	Information about resource related location information	string string anyURI float + float	EBUCoreMain/coreMetadata/coverage/dc:coverage EBUCoreMain/coreMetadata/coverage/spatial/location/name EBUCoreMain/coreMetadata/coverage/spatial/location/code EBUCoreMain/coreMetadata/coverage/spatial/location/posx + EBUCoreMain/coreMetadata/coverage/spatial/location/posy
ma:compression	encoding/@typeLabel encoding/@typeLink	free text or a link to a classification scheme e.g. published as a web resource	string anyURI	EBUCoreMain/coreMetadata/format/channel/encoding/@typeLabel EBUCoreMain/coreMetadata/format/channel/encoding/@typeLink
ma:frameSize	height width	see the syntax of ma:frameSize for correct mapping	nonNegativeInteger nonNegativeInteger	EBUCoreMain/coreMetadata/format/height EBUCoreMain/coreMetadata/format/width

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

MAWG	EBU Core	How to do the mapping	Datatype	XPath
Note: each is to be preceded by the namespace 'ebucore:'				
ma:targetAudience	dc:type genre/@type Label genre/@type Link	Free text in type or genre/@typeLabel or @typeLink (using a targetAudience Classification Scheme or equivalent)	string string anyURI	EBUCoreMain/coreMetadata/type/dc:type EBUCoreMain/coreMetadata/type/genre/@typeLabel EBUCoreMain/coreMetadata/type/genre/@typeLink
ma:locator	Format/Location	an address at which the resource can be found and e.g. played from e.g. a dvb url	string	EBUCoreMain/coreMetadata/location
ma:frameRate	frameRate	if in dc:format, a syntax should be used to prefix the property being documented e.g. frameRate:xxx	string	EBUCoreMain/coreMetadata/format/dc:format
ma:samplingRate	samplingRate	if in dc:format, a syntax should be used to prefix the property being documented e.g. samplingRate:xxx	string	EBUCoreMain/coreMetadata/format/dc:format
ma:bitrate	bitrate	if in dc:format, a syntax should be used to prefix the property being documented e.g. bitrate:xxx	string	EBUCoreMain/coreMetadata/format/dc:format
ma:numTracks	videoFormat audioFormat	the video or audio formats imply the number of video and / or audio tracks	string anyURI string anyURI	EBUCoreMain/coreMetadata/format/videoFormat/@formatLabel EBUCoreMain/coreMetadata/format/videoFormat/@formatLink EBUCoreMain/coreMetadata/format/audioFormat/@formatLabel EBUCoreMain/coreMetadata/format/audioFormat/@formatLink
ma:fragment	hasPart	e.g. a scene or shot identified by its uri	string anyURI	EBUCoreMain/coreMetadata/hasPart/relationLink
ma:namedFragments	hasPart	e.g. a scene or shot identified by an identifier or title	string string	EBUCoreMain/coreMetadata/hasPart/dc:relation EBUCoreMain/coreMetadata/hasPart/relationIdentifier/dc:identifier

Table 4: Mapping between EBU Core and W3C Ontology for Media Resource (MAWG)

TV Anytime

TV-Anytime (TVA) Metadata has been designed to support the Business-to-Consumer exchange in the broadcast industry [TVA1,TVA2]. It allows the consumer to find, navigate and manage content from a variety of internal and external sources including, for example, enhanced broadcast, interactive TV, Internet and local storage. Metadata is generated during the process of content creation and content delivery. There are three basic kinds of metadata: *Content Description*, *Instance Description*, and *Consumer Metadata*. In addition the standard defines Segmentation Metadata and Metadata Origination Information Metadata. The information that the consumer or agent will use to decide whether or not to acquire a particular piece of content is called *attractors*, and is used in electronic programme guides, or in Web pages. These attractors rely on descriptors stemming from MPEG-7. Furthermore, some MPEG-7 data types are used directly (e.g., mpeg7:TextualType is used for many TVA of elements.)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

The *content description metadata* describes content independently of any particular instantiation of a media programme. Programme in this context means an editorially coherent piece of content.

Descriptions of content, e.g., television programmes are held in the *ProgramInformationTable*. They include metadata like the Title (here the mpeg7:TitleType is used) of the programme, a Synopsis, the Genre it belongs to, and a list of Keywords that can be used to match a search. Descriptions of groups of related items of content e.g. all episodes of "Foxes in the Wild" are held in the *GroupInformationTable*. They include among other the GroupType, a BasicDescription and MemberOf element. A mapping of cast members to unique identifiers is held in the *CreditsInformationTable*. The identifiers can be used in other metadata instances simplifying the search. The purchase information, like Price and PurchaseType, is held in the *PurchaseInformationTable*. Critical reviews of items of content are held in the *ProgramReviewTable*. They include metadata like the Reviewer, a FreeTextReview and the ProgramId.

Instance Description Metadata is required in case of significant differences between instantiations of the same content. These are instances with the same CRID (the CRID connects content metadata with content). Instance Metadata is connected with content related to a definite event. Descriptions of particular instances (locations) of content are held in the *ProgramLocationTable*. They include the elements Schedule, BroadcastEvent, OnDemandProgram and OnDemandService, all derived from ProgramLocationType. They include among other information about the programme, start and end. Also Title, Synopsis, Genre and PurchaseList can be specified. Descriptions of services within a system are held in the *ServiceInformationTable*. For each single Service Name, Owner Logo (here the mpeg7:MediaLocatorType is used), ServiceDescription, ServiceGenre etc. can be specified.

Consumer Metadata includes Usage History and User Preferences, both based on respective MPEG-7 data types. The Usage History provides a list of the actions carried out by the user over an observation period. It is used for tracking and monitoring the content viewed by individual members. Thus, it builds a personalized TV guide by tracking user viewing habits, selling viewing history to advertisers or tracking and monitoring content usage for more efficient content development. The User Preferences facilitate description of user's preferences pertaining to consumption of multimedia material. They include FilteringAndSearchPreferences and BrowsingPreferences and can be correlated with media descriptions to search, filter, select and consume desired content.

VideoActive

In an earlier section we discussed extensively the Europeana Data Model EDM (see p. 15). The VideoActive project⁴⁴ has shown how descriptive metadata of TV programs can be modelled with EDM. The figure below shows an example. This example could also be used as a model for the way in which PrestoPRIME metadata could be exported to Europeana. The figure shows that provenance metadata are attached to the "aggregation" and descriptive metadata are attached to the "Proxy". The metadata elements are based on the DCTERMS elements set.

⁴⁴ <http://videoactive.wordpress.com/> (verified 24 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

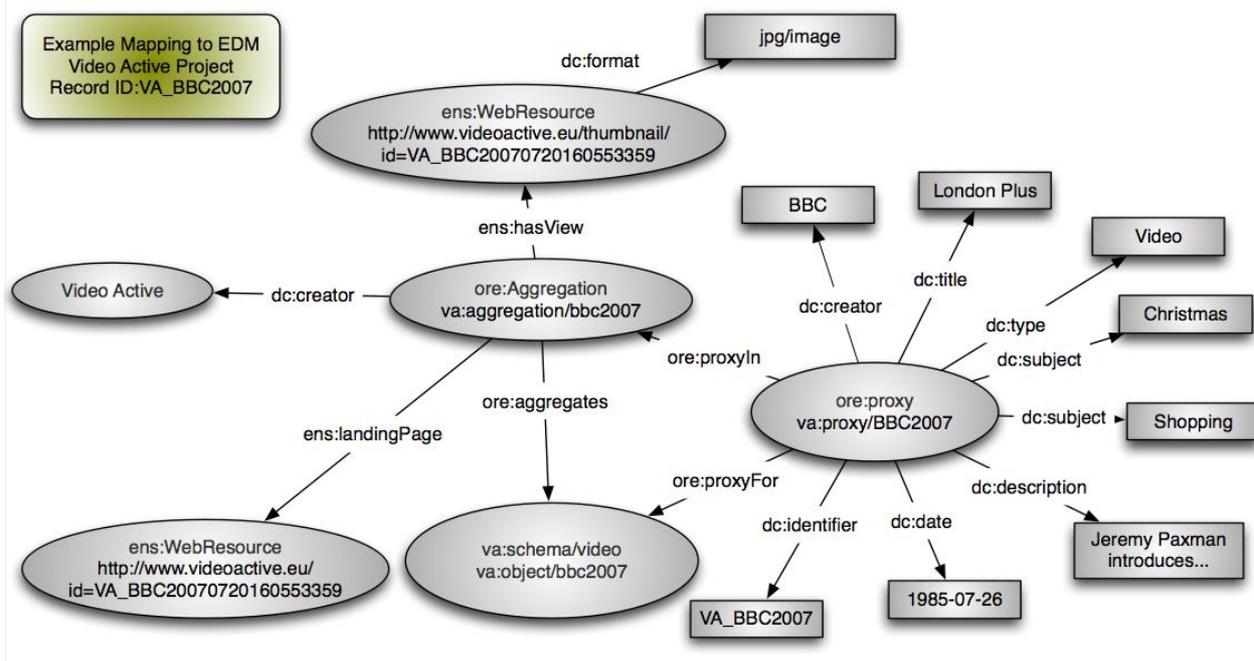


Figure 4: VideoActive example of EDM usage (provided by Vassilis Tzouvaras, National Technical University of Athens, copied with permission)

VRA

VRA Core 4.0⁴⁵ is a data standard for the cultural heritage community that was developed by the Visual Resources Association's Data Standards Committee. It consists of a metadata element set (units of information such as title, location, date etc.), as well as an initial blueprint for how those elements can be hierarchically structured. The element set provides a categorical organisation for the description of works of visual culture as well as the images that document them.

The VRA 4.0 Core Categories are defined as a specialization of the DCMI terms (see page 9), using the Dublin Core dumb-down principle. This means that VRA can be used interoperable by an organisation that has adopted the DCMI Terms standard.

VRA provides an explicit distinction between metadata about a “work“ on the one hand and an image (of the work) on the other hand. This distinction is not explicit in Dublin Core, but it is in line with the EDM model, where (following OAI ORE) objects are distinguished from their digital representations.

PBCore

According to the Public Broadcasting Metadata Dictionary (PBCore) website⁴⁶ PBCore intends to be:

- a core set of terms and descriptors (elements)...
- used to create information (metadata)...
- that categorises or describes...
- media items (sometimes called assets or resources).

⁴⁵ <http://www.vrweb.org/projects/vracore4/index.html> (verified 2 June 2010)

⁴⁶ <http://www.pbcore.org/PBCore/index.html> (verified: 10 May 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

Within a project – the PBMD Project – lasting over two years various standards, dictionaries and schemes for metadata description were compared. Aim was to arrive at the smallest set of descriptors or elements that could adequately describe and catalogue media items which are produced at Public Broadcasting radio and television stations. Information should be sharable between stations, regional and national distributors, independent producers and also vendors of Digital Asset Management systems. In 2005 version 1.0 was launched and an update – version 1.1 – followed in 2007.

The basis of PBCore is Dublin Core (ISO 15836). The work has also been reviewed by the Dublin Core Metadata Initiative Usage Board. PBCore includes 53 elements and they are arranged in 15 containers and 3 sub containers. Descriptions in PBCore allow 4 content classes to be described as follows:

- **PBCoreIntellectualContent** (9 containers; 16 elements) for describing the actual intellectual content of a media asset or resource
- **PBCoreIntellectualProperty** (4 containers; 7 elements) related to the creation, creators, usage, permissions, constraints and use obligations associated with a media asset or resource
- **PBCoreInstantiation** (1 container; 3 sub-containers; 28 elements) for identifying the nature of the media asset (existing in some form or format in the physical world or digitally)
- **PBCoreExtensions** (1 container; 2 elements) cover additional descriptions that have been crafted by organisations outside of the PBCore Project. These extensions fulfil the metadata requirements for these outside groups as they identify and describe their own types of media with specialised, custom terminologies unique to their needs and community requirements.

There are currently no elements in PBCore which allow structuring of audiovisual content. PBCore is undergoing a re-design to become version 2.0⁴⁷.

2.3 Models for preservation metadata

According to the PREMIS data dictionary [PREMIS2008] five entities can be defined which are said to be important for modelling information on digital preservation activities: the entities are Intellectual Entities, Objects, Events, Rights, and Agents. So these five entities could be a starting point to think about the metadata elements to be listed and used.

Major questions at the beginning of setting up a concrete model for preservation metadata would be:

- How will the archives use the preservation metadata?
- Will they try to come to some common understanding of what their preservation metadata should look like?
- Do they already have some kind of preservation metadata, and which format is used for that purpose?
- What are the expectations for standardisation of a preservation metadata representation?
- For which scenarios defined in PrestoPRIME's WP5 will the preservation metadata be relevant and how can they be applied therein?

⁴⁷ <http://pbcore.org/2.0/> (verified: 11 May 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

PREMIS

The following subsections show the principal entities as defined in the PREMIS data dictionary with the top level attributes foreseen therein.

Intellectual Entities

The Intellectual Entity is considered out of scope for PrestoPRIME because it is well served by descriptive metadata.

Objects

The Object entity aggregates information about a digital object held by a preservation repository and describes those characteristics relevant to preservation management.

- 1.1 objectIdentifier (M⁴⁸, R⁴⁹)
- 1.2 objectCategory (M, NR)
- 1.3 preservationLevel (O, R) [representation, file]
- 1.4 significantProperties (O, R)
- 1.5 objectCharacteristics (M, R) [file, bitstream]
- 1.6 originalName (O, NR) [representation, file]
- 1.7 storage (M, R) [file, bitstream]
- 1.8 environment (O, R)
- 1.9 signatureInformation (O, R) [file, bitstream]
- 1.10 relationship (O, R)
- 1.11 linkingEventIdentifier (O, R)
- 1.12 linkingIntellectualEntityIdentifier (O, R)
- 1.13 linkingRightsStatementIdentifier (O, R)

Events

The Event entity aggregates information about an action that involves one or more Object entities. Metadata about an Event would normally be recorded and stored separately from the digital object.

Whether a preservation repository records an Event or not depends upon the importance of the event. Actions that modify objects should always be recorded. Other actions such as copying an object for backup purposes may be recorded in system logs or an audit trail but not necessarily in an Event entity.

- 2.1 eventIdentifier (M, NR)
- 2.2 eventType (M, NR)
- 2.3 eventDateTime (M, NR)
- 2.4 eventDetail (O, NR)
- 2.5 eventOutcomeInformation (O, R)

⁴⁸ Obligation: **Mandatory** or **Optional**

⁴⁹ Repeatability: **Repeatable** or **NotRepeatable**

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

2.6 linkingAgentIdentifier (O, R)

2.7 linkingObjectIdentifier (O, R)

2.7.1 linkingObjectIdentifierType (M, NR)

2.7.2 linkingObjectIdentifierValue (M, NR)

2.7.3 linkingObjectRole (O, R)

Rights

For the purpose of the PREMIS Data Dictionary, statements of rights and permissions are taken to be constructs that can be described as the Rights entity. Rights are entitlements allowed to agents by copyright or other intellectual property law. Permissions are powers or privileges granted by agreement between a rightsholder and another party or parties.

A repository might wish to record a variety of rights information including abstract rights statements and statements of permissions that apply to external agents and to objects not held within the repository. The minimum core rights information that a preservation repository must know, however, is what rights or permissions a repository has to carry out actions related to objects within the repository. These may be granted by copyright law, by statute, or by a license agreement with the rightsholder.

4.1 rightsStatement (O, R)

4.2 rightsExtension (O, R)

Agents

The Agent entity aggregates information about attributes or characteristics of agents (persons, organisations, or software) associated with rights management and preservation events in the life of a data object. Agent information serves to identify an agent unambiguously from all other Agent entities.

3.1 agentIdentifier (R, M)

3.2 agentName (O, R)

3.3 agentType (O, NR)

METS (wrapper format)

Whereas in general “XML has become the de-facto standard for representing metadata descriptions of resources on the Internet” (from [HUNTER]) a standard representation for expressing the hierarchical structure of digital library objects is the Metadata Encoding and Transmission Standard (METS)⁵⁰. Therefore METS is not a standard for preservation metadata but a wrapper format, i.e. an XML Schema designed for purpose of including the names and locations of the files that comprise objects and their associated metadata [CUNDIFF].

The main characteristics are: open standard, non-proprietary, developed by the library community, (relatively) simple, modular and extensible.

The development of METS' predecessor (called MOA2) was started in 1997 with the *Making of America II* initiative. The goal of MOA2 was to create a digital object standard for encoding structural, descriptive and administrative metadata along with the primary content. Later additional needs emerged (e.g. support for time-based content and more

⁵⁰ <http://www.loc.gov/standards/mets/> (checked: 25 June 2009)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

flexibility in descriptive and administrative metadata). The previous MOA2 XML DTD was revised and resulted in the METS XML schema version 1.2 in 2002. Its current version is 1.8 from April 2009.

The Metadata Encoding and Transmission Standard (METS) allows the use of externally developed metadata schemes. They can be fit into its two defined metadata sections, <dmdSec> and <amdSec>. METS itself does not care about the descriptive or administrative metadata schemes that are incorporated by implementers.

Some community based standards are recognized by the METS board and PREMIS is one of them. It should be used as administrative metadata within METS. To use PREMIS together with METS some decisions have to be made in advance as the PREMIS schema was developed in an implementation neutral way and also METS has quite some flexibility within it.

Several issues are to be thought about:

- there exist redundancies within PREMIS and METS
- PREMIS elements can be used in a number of METS sections (a related issue is the question of how many METS sections are used)
- whether to use the PREMIS container schema or not
- how to deal with format specific metadata within PREMIS

As a consequence on these issues guidelines have been developed between the METS and PREMIS communities. The working draft is available for experimentation and comment at [PREMISMETS2008].

Metadata model of the New Zealand National Library (based on PREMIS and METS)

The National Library of New Zealand (NLNZ) has been a pioneer in setting up approaches to the management of electronic material and hence the need for a Digital Archive and further the desire to improve access to its collections via digitisation.

The need for preservation of digital materials became also obvious. Hand in hand with that the business need for preservation metadata emerged as an integral component.

The Digital Archive at NLNZ enhances access to the Library's digital resources for all New Zealanders which is necessary if the Library is to achieve its mandate 'to collect, preserve and make available recorded knowledge, particularly that relating to New Zealand,' in an environment increasingly characterised by electronic online and offline resources.

The NLNZ preservation metadata schema details the data elements needed to support the preservation of digital objects and forms the basis for the design of a database repository and input systems for collecting and storing preservation metadata. It incorporates a number of data elements needed to manage the metadata in addition to metadata relating to the digital object itself. The aim has been to produce a document that will serve as an implementation template while at the same time remaining consistent with standards being developed internationally around preservation metadata. (from [NLNZ2003])

The preservation metadata at NLNZ will be used to:

- store information supporting preservation decisions and actions
- document preservation processes, such as migrations, transformations and emulations

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

- record the effects of preservation processes
- ensure the authenticity of Preservation Masters over time
- enable objects for which the library has assumed preservation responsibility to be identified.

And further Preservation metadata addresses two functional objectives:

1. providing the Library with sufficient knowledge to take appropriate actions in order to maintain a digital object's bit stream over the long-term.
2. ensuring that the content of an archived object can be rendered and interpreted, in spite of future changes in storage and access technologies.

The NLNZ Preservation Metadata Set

The library set up a Preservation Metadata Set and believed that it should include information which was needed to preserve the digital collections. The model indicated which information will be needed and used in the future but not what data were going into the system and how / when / by whom this could be done. Neither did the model point out how the metadata were associated to the described objects. The model should be open to a variety of applications recording that information. So the definition was mainly outcome driven and the model simply said: *"however you do it, this is what you have to deliver so we can manage preservation"*.⁵¹ This was before the year 1999.

Referring to [THOMPSON2003] the resulting schema at NLNZ took into account three aspects:

1. **limiting the scope of preservation metadata** to only those data required for digital preservation. Elements supporting other activities like preservation of analogue formats or resource discovery were stripped out. Elements common between preservation and other functions were identified for repository element. Further removed was the need to collect preservation metadata about dissemination formats.
2. **maximising potential for automation** (i.e. automatic population of the maximum number of elements) is demonstrated with the focus on the Preservation Master (which is not the "original" but the "best effort" representation of the material). The Preservation Master is the subject to preservation processes being transformed from obsolete formats into current ones. The use of preferred file types allows more standardised processes and their number becomes limited. With this approach the range of values collected as preservation metadata will be reduced and will thus be more manageable sets referring to automation.
3. **ensuring change control for metadata** was implemented as well. The implemented audit trail includes information which allows tracking modifications of metadata along with information about the persons who were responsible for the changes.

National Digital Heritage Archive (NDHA)

Later in May 2004, the National Library of New Zealand was allocated NZ\$24 million [US\$16 million] by the national government to fund a program to establish a trusted digital repository to protect the nation's digital documentary heritage for future generations. The

⁵¹ Preservation Metadata for Digital Collections; <http://www.nla.gov.au/preserve/pmeta.html> (checked: 21 July 2009)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

new to establish National Digital Heritage Archive (NDHA) was to collect, preserve and make accessible digital objects both online and offline, including websites, published works, images and material contained on CDs and floppy disks.⁵²

So digital preservation at the Library of New Zealand was not seen as a once-only activity. The aim was to preserve the heritage collections “in perpetuity,” i.e. a really long time. Further this means that the systems developed and installed were to manage and preserve digital material over time and therefore must be capable to evolve over time including access to all the digital objects in the future. To do so the programs have to be robust and should align with international standards and best practices in that area. The programs have to be open to new formats or changing preservation strategies. Also they had to be adoptable to various organizations.

To achieve this, a range of activities, emerging standards and best practice initiatives were taken into account for the development of the NDHA. They include:

- Web archiving tools — IIPC, Nordic Web Archive and PANDORA. Together with the British Library the so-called Web Curator Tool⁵³ was developed.
- Preservation metadata — NLNZ and PREMIS.
- Structural metadata — METS/MPEG 21.
- Persistent identifiers — Handle/DOI.
- Rights management — INDECS.
- File format identification, metadata extraction — NLNZ and JHOVE.
- Digital preservation R&D — CAMiLEON, Digitale Duurzaamheid, NDIIPP, Xena and Variable Media Network.

The developments around NDHA were done by NLNZ in a partnership with the Ex Libris Group and with Sun Microsystems. The first phase of the NDHA was finished in October 2008 and the final phase will be completed by 2010.⁵⁴ Further information on the project can be found in the case study⁵⁵. It provides information on the initial phase (collection of requirements), the way of cooperating (buy, build, or being partner), how to build the solution from its parts (for ingest, storage, data management, administration, preservation planning and access) and architecture (hardware and software) itself. One outcome of the project was the Rosetta system from Ex Libris previously known as the Digital Preservation System (DPS).⁵⁶ Rosetta was also chosen by the Bavarian State Library as a Strategic Partnership with Ex Libris for Long Term Preservation in November 2009⁵⁷. Some reasons for that were that the system is based on OAIS and conforming to trusted

⁵² Steve Knight, Manager Innovation Centre, Digital Innovation Services, National Library of New Zealand, Wellington, New Zealand In Perpetuity: A Nation's Well-Spring of Knowledge; <http://www.elsevier.com/wps/find/librariansinfo.librarians/LCN030404> (checked: 21 July 2009)

⁵³ Web Curator Tool <http://www.natlib.govt.nz/services/get-advice/digital-libraries/web-curator-tool> (checked: 14 December 2009)

⁵⁴ National Digital Heritage Archive <http://www.natlib.govt.nz/about-us/current-initiatives/ndha> (checked: 19 October 2009)

⁵⁵ Case Study: Digital Preservation at the National Library of New Zealand (Preservation: A Forward-Looking Mission) <http://www.exlibrisgroup.com/files/CaseStudy/SunPreservationandNLNZ.pdf> (checked: 14 December 2009)

⁵⁶ Press Release: Ex Libris Group Announces the General Release of its Digital Preservation System <http://www.exlibrisgroup.com/?catid=%7b916AFF5B-CA4A-48FD-AD54-9AD2ADADEB88%7d&itemid=%7b9B1F2C8A-3B03-459F-A2B4-4425A4D79689%7d> (checked: 15 July 2009)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

digital repository (TDR) requirements. Other crucial factors were the functionality for preservation planning module and compliance with international preservation and digitization standards such as METS.

DNX

Several approaches for implementing PREMIS were presented at the **PREMIS Implementation Fair** (sponsored by the Library of Congress and held on October 7, 2009).⁵⁸ One of those approaches is the use of PREMIS in Rosetta (from Ex Libris).

The data model of Rosetta is based on PREMIS and the system is compliant both with PREMIS and METS. METS documents in Rosetta describe intellectual entities consisting of several representations. Part of the data model of Rosetta is DNX a proprietary schema from Ex Libris for recording technical metadata. DNX includes and normalises format-specific technical metadata from multiple schemas like MIX (NISO Metadata for Images in XML)⁵⁹, textMD (Technical MetaData for text)⁶⁰. Further it can include PREMIS information. Some redundancies exist between DNX, METS, and PREMIS. As of late Autumn 2009 Rosetta does not have the ability to convert from a DNX file to a PREMIS record but this functionality is foreseen at Ex Libris for the future in the case when there is an established PREMIS community who would like to make use of it.

Long-term preservation Metadata for Electronic Resources (LMER)

In parallel to early PREMIS developments by OCLC and RLG, the German National Library started an initiative to work on metadata for long-term preservation. The initiative was called "*Langzeitarchivierungsmetadaten für elektronische Ressourcen*" (which means "*Long-term preservation Metadata for Electronic Resources*").

The trigger was as in other cases missing or insufficient possibilities and standards for collecting technical metadata which can be used for long-term preservation of electronic documents. Therefore the German National Library started with the new schema – LMER. The schema can be found along with other information on the website of LMER⁶¹. The base for this schema was the model which was used in the National Library of New Zealand. As a consequence the schema included sections like in that metadata model. Further the overlapping PREMIS activities were seen as relevant for the work on the schema and the wrapping mechanism of METS introduced.

Shared format registries

Preservation of digital objects deals a lot with file formats and this also means rather complex information (format, version, software versions for support etc.). To use them in a consistent way in different archives but even within one organization a central lookup (i.e. at a registry) is the best way. Using the service of registry an individual repository just

⁵⁷ Press Release: The Bavarian State Library Selects the Rosetta Digital Preservation System and Establishes a Strategic Partnership with Ex Libris for Long Term Preservation http://www.exlibrisgroup.com/default.asp?catid=%7b916AFF5B-CA4A-48FD-AD54-9AD2ADADEB88%7d&details_type=1&itemid=%7b9F6528E1-DE3B-447E-B431-B73B43E271B7%7d (checked: 14 December 2009)

⁵⁸ <http://www.loc.gov/standards/premis/pif-presentations/PREMISImplementationFairSummary.pdf> (checked: 11 May 2010)

⁵⁹ <http://www.loc.gov/standards/mix/> (checked: 11 May 2010)

⁶⁰ <http://www.loc.gov/standards/textMD/> (checked: 11 May 2010)

⁶¹ <http://www.d-nb.de/standards/lmer/lmer.htm> (checked: 1 July 2009)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

stores information on the file format and where/how to find the detailed information referred to. A number of such registries exist or are developed.

One of these registries is the online service PRONOM⁶² hosted by the UK National Archives. The registry stores information about file formats, software products, software vendors, and product support providers. The service is currently designed for human-interactive use but not for being queried by programs. When asked for a particular file format PRONOM lists detailed technical metadata accompanied by a list of software applications that can render the requested format. Search functionality allows queries for formats, products, vendors and by the PRONOM identifier. By using the specific "lifecycle" search all software products supported as of a certain date, or released before or after a given date can be found.

Up to now most information in the database was collected by preservation staff at the National Archives. For the future the hope is that developers of software products and file format specifications submit information on new formats, versions etc. directly to PRONOM. Along with the registry information The National Archives intend to provide tools like e.g. DROID (Digital Record Object Identification), a free tool for identifying file formats.

The Global Digital Format Registry (GDFR)⁶³ was another registry with similar purpose, i.e. providing information on file formats and so on.

In April 2009 the GDFR initiative joined forces with the UK National Archives' PRONOM registry initiative under a new name – the Unified Digital Formats Registry (UDFR)⁶⁴. The aim was to deal with enlarged requirements of a larger digital preservation community. The UDFR will support the requirements⁶⁵ and use cases compiled for GDFR and will be seeded with PRONOM's software and formats database.

2.4 Models for provenance metadata

Most of the literature in the field comes from the area of e-science. With the wider adoption of Linked Open Data⁶⁶, provenance information becomes a crucial issue on the Semantic Web. The W3C has thus chartered the Provenance Incubator Group (ProvXG)⁶⁷ to survey the state of the art and define the requirements for handling provenance information on the Web. The incubator group has also collected references to provenance related literature and tools which cover a much wider scope than discussed here.

This section focuses on representation of provenance information from a metadata model perspective, and not on automatic creation of provenance information. However, such approaches (as e.g. proposed in [FMS07,GM09]) might also be relevant inside the PrestoPRIME preservation system.

⁶² PRONOM registry <http://www.nationalarchives.gov.uk/pronom/> (checked 17 July 2009)

⁶³ Global Digital Format Registry Information Site <http://www.gdfr.info/> (checked 22 July 2009)

⁶⁴ Unified Digital Formats Registry <http://www.udfr.org/> (checked 22 July 2009)

⁶⁵ http://gdfr.info/wiki/index.php/Activity_3:_Gathering_of_functional_and_non-functional_requirements_for_a_community_format_registry (checked 11 June 2010)

⁶⁶ <http://linkeddata.org/> (checked 21 June 2010)

⁶⁷ <http://www.w3.org/2005/Incubator/prov> (checked 21 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

ORE

The ORE data model⁶⁸ can serve to represent basic provenance information of a compound resource. ORE defines aggregations that express a grouping of resources. A Resource Map describes an aggregation and can hold for example metadata about the creation of one or several resources in the map (e.g. dc:creator, dcterms:created, pointers to rights description). Aggregations can be nested into arbitrary graphs.

PREMIS

Although not being its primary function, PREMIS [PREMIS2008] does provide some support for provenance metadata. The typical approach (e.g. described in [Hab07,Hoe09]) is to use PREMIS events and agents to represent content creation and modification events. This of course requires ideally some controlled vocabulary for the types of events needed in provenance information.

The PREMIS implementation report [Woo07] states that the New Zealand National Digital Heritage Archive also uses a trail of events to document provenance information, but just describing a few properties per event. This report also points out that PREMIS supports description of derivation relationships (in contrast to just structural relationships as in METS), which can also be used to express provenance of content items.

In [Hoe09] an alternative approach is described, which proposes to use the extension elements of *Object's objectCharacteristics* or *significantProperties* to include elements from a domain specific provenance vocabulary.

METS

METS⁶⁹ has a specific element for provenance metadata in the administrative metadata section, called *digiProvMD*, intended to describe any preservation-related actions on the digital object. It is a generic element that can wrap or reference metadata describing these actions. Although it is sufficiently generic to embed any metadata format, it is due to its definition not appropriate for describing provenance information of the original resource (whether born analogue or digital).

In the context of the LOC's Digital Audio-Visual Preservation Prototyping Projects⁷⁰ an extension schema called PMD⁷¹ for describing preservation related actions of a/v objects has been proposed. The schema defines the concept hierarchy Process – Task – Tools – Settings to describe actions on the digital object.

Open Provenance Model (OPM)

The OPM⁷² [Mor07] aims at defining an interoperable, technology-agnostic provenance model for exchange between systems. The authors specify the model, but not any specific representation or serialisation. The three main entities are Artifact, Process and Agent. In addition the model defines six relations and roles. Optionally, time can be annotated and used for inference, e.g. the time of the result of a causal relation must be larger than the time of the cause. XML and OWL representations of OPM have been defined.

⁶⁸ <http://www.openarchives.org/ore/1.0/datamodel> (checked 17 June 2010)

⁶⁹ <http://www.loc.gov/standards/mets/> (checked 17 June 2010)

⁷⁰ <http://www.loc.gov/rr/mopic/avprot/avprhome.html> (checked 17 June 2010)

⁷¹ http://www.loc.gov/rr/mopic/avprot/digiprov_expl.html (checked 17 June 2010)

⁷² <http://openprovenance.org/> (checked 17 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

Changeset

Changeset⁷³ is a small RDF vocabulary for describing the change history of resources. Each Changeset captures date, author, reason and subject of the change. Changesets can be linked to form a history.

Provenance Vocabulary

The Provenance Vocabulary⁷⁴ is an RDF vocabulary for describing provenance information. The three main entities of the model are Artifact, Actor and Execution. Several specialisations of these entities are defined, and their relations can be amended with the description of data, guidelines etc. used in an execution.

Provenir

The Provenir Ontology⁷⁵ is a provenance vocabulary targeting the e-science domain. The ontology defines three main classes (data, process and agent), five specialisations of these classes and nine properties relating these classes.

Content identifiers

A topic related to the handling of provenance information is the unique identification of audiovisual content.

ISAN

The ISAN (International Standard Audiovisual Number) is only applicable to complete commercial works. It is not applicable to user-generated contents or excerpts, and can be lost when successive copies are made.

IMDb (Internet Movie Database)

IMDb is a user-contributed registry of one million programmes (50% television programmes, 50% movies), but with no reference to contents, and no link to ISAN; a considerable mass of content is available for viewing on YouTube, but there is no reference to a registry.

2.5 Models for rights metadata

Rights metadata are handled intensively in another task of PrestoPRIME (WP4T4). The corresponding report is PrestoPRIME deliverable D4.0.5.

⁷³ <http://vocab.org/changeset/schema.html> (checked 17 June 2010)

⁷⁴ http://sourceforge.net/apps/mediawiki/trdf/index.php?title=Provenance_Vocabulary (checked 17 June 2010)

⁷⁵ http://wiki.knoesis.org/index.php/Provenir_Ontology (checked 17 June 2010)

3 Shortcomings and gaps in metadata standards and models

A variety of different standards or data models exist for various purposes: descriptive metadata, preservation metadata, provenance metadata and rights metadata.

Descriptive metadata are mainly dealt with in local systems – either standardised or in proprietary formats. In the previous PrestoSpace project a data model was defined including MPEG-7 and P_Meta. Further representations are in use in the archives but they can be mapped to the aforementioned model.

As previously mentioned some standards exist for describing preservation, provenance and rights. An obvious candidate – developed for those purposes – is PREMIS. That standard provides four entities: objects, events, agents and rights. Wrapping such information with other information can be done in METS.

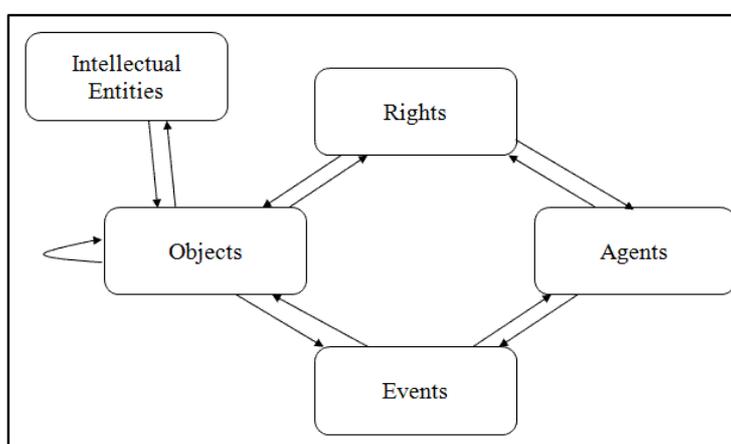


Figure 5: the PREMIS entities

The PREMIS Metadata Dictionary [PREMIS2008] provides information on the possibilities when using the PREMIS standard for preservation relevant information. The section on “Limits to the scope of the Data Dictionary” tells in detail what is not possible to do within particular domains.

3.1 Preservation metadata

Mostly information for general preservation activities can be stored in PREMIS. The PREMIS working group did not dig into the details of very different domains but stayed at a level which can be used independent from the domain where it should be used in.

The coverage of technical metadata as would be necessary for a/v content was not an aim in the working group due to a lack of time and also due to missing expertise in the working group. The description of media and hardware cannot be documented in detail in PREMIS. A possibility in PREMIS would be to use the extension with the semantic unit *objectCharacteristicsExtension*. It allows for including information given in an external scheme with a deeper granularity. Technical metadata could be represented in EBU tech 3295 (P_Meta⁷⁶) or SMPTE RP210⁷⁷. These schemes could be included as mentioned before. In practice, it can be expected that properties from more than one such vocabulary are needed. Neither METS nor PREMIS define a generic way of including metadata properties except for including XML elements from another schema. Both P_Meta and SMPTE RP210 define their own XML representations, i.e. the P_Meta XML schema and

⁷⁶ <http://tech.ebu.ch/docs/tech/tech3295v2.pdf> (checked 24 June 2010)

⁷⁷ <http://www.smpete-ra.org/mdd/index.html> (checked 24 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

SMPTE 434M [MXFXML] for representing MXF header structures in XML. The key/value structure of DNX might be an option for representing technical properties stemming from different metadata vocabularies.

Regarding the modelling of preservation options, the inclusion of business rules in PREMIS was mentioned. This could be done with an entity similar to the Rights entity but such an entity has not been defined to date. One exception exists in PREMIS: the *preservationLevel* (Information indicating on the set of preservation functions to be applied to an object) within the Object entity was felt as important.

3.2 Provenance metadata

All of the more comprehensive provenance metadata models have a representation of the three entities: object, event and agent (although they might be named differently). Some of these models allow these entities to have a certain set of properties, but specialisations can be designed at least for those models that are based on RDF. The models differ in the set of properties between these entities, but – again – defining sub properties is possible in the RDF based models if necessary. Some of the models are defined as abstract models or vocabularies, not bound to a specific representation, so that an appropriate representation (i.e. one that can be serialised in XML) needs to be selected.

The critical issues w.r.t. provenance metadata models are the following:

- In media production and preservation processes we encounter a wide range of actions/events/processes involving complex tools with many settings. The format to be used must be able to define a taxonomy of actions/events/processes (either by using type attributes or sub classing) and must allow representing (some of) their parameters.
- The objects involved in provenance-related events are often small fragments of media items. The model must thus be able to reference arbitrary fragments of media items, ideally using a universal identifier such as a Media Fragment URI⁷⁸.

The implementation of provenance information for a/v content within the possibilities of PREMIS is more or less limited to the relationship and linking to other objects. Relationships in PREMIS can be expressed as structural or as derivation relationships. Structural is meant for defining parts of objects. Derivation may be the result of replication or transformation of an object. Both these relationships could be used for a/v content. PREMIS supports links to objects on the file or bitstream level; however, neither seems to be sufficient for referencing meaningful segments. It could be an option to use Media Fragment URIs⁷⁹ as types of `objectIdentifierType` in order to reference objects that represents temporal/track fragments of audiovisual media. An alternative would be to use fragment identifiers define in the descriptive metadata, however, this introduces a dependency between provenance and descriptive metadata.

Another issue with PREMIS is the description of events. `eventType` uses a controlled vocabulary that could be adapted to the needs of PrestoPRIME. However, for many events, tool parameters are important to be documented. PREMIS `eventDetail` is intended to be only an informative, non-machine processable description of event parameters. In order to use it, a machine processable description of event parameters is needed.

⁷⁸ <http://www.w3.org/TR/media-frags/> (checked 17 June 2010)

⁷⁹ <http://www.w3.org/2008/WebVideo/Fragments/> (checked 17 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

A third issue is the embedding of provenance metadata into METS. According to the METS schema specification, `digiprovDM` is defined to “*record any preservation-related actions taken on the various files which comprise a digital object (e.g., those subsequent to the initial digitization of the files such as transformation or migrations) or, in the case of born digital materials, the files’ creation*”. This introduces an inconsistency, as provenance information of born digital content is entirely covered, while provenance information related to analogue production (or analogue preservation actions) would have to be represented in `sourceMD`.

3.3 Rights metadata

Rights metadata are dealt with in PrestoPRIME task WP4T4. Deliverable D4.0.5 reports about gaps w.r.t. rights metadata and proposes appropriate solutions.

4 Conclusions and Recommendations

4.1 *Integration of preservation and provenance metadata models*

Several models allow storing the preservation related metadata, provenance information and rights information. For covering all necessary pieces of information some of the models have to be combined in some sense.

The different representations for information on preservation (including technical metadata), on provenance and on rights should be placed within one containing format. A commonly used format for that purpose is METS. This is due to several of its built-in features but also due to its ability to host other formats as long as they are described in an XML schema.

The various aspects have to be covered in the combined model. In the METS container the information is included either as an included XML block or as a reference to an external information block.

- Content identification: METS provides only onboard means for identifying the container, for content identifiers either DCTerms:identifier or the objectIdentifier of the PREMIS object entity in `amdsec/sourceMD` section of METS can be used
- a/v content description is done with different representations and stored in METS in `dmdsec` either as included information or a reference to external information. Different representations can be mapped with appropriate mechanisms
- preservation metadata: general information is covered in PREMIS, and PREMIS metadata can be incorporated within a METS container. The METS guidelines⁸⁰ for integration of PREMIS metadata in METS recommend either including the whole PREMIS block under `amdsec/digiprovMD` or including the PREMIS first level data elements in appropriate sections (also listed)
- technical metadata (P_Meta or SMPTE RP210 integrated in the `amdsec/techMD` section in METS)
- preservation options: PREMIS
- provenance (integrated in the `amdsec/digiprovMD` section of METS)
- rights (are handled with MPEG-21 and [MVCO] extensions and can be placed into `amdsec/rightsMD` section of METS)

4.2 *Relation to Europeana*

Europeana was identified by PrestoPRIME as a primary dissemination channel for the project. Therefore, we have paid special attention to interoperability issues in connection with Europeana.

TV archives that want to expose their content through Europeana will want to have a clear route into the Europeana ingestion process. Of course, not all metadata of objects in an a/v archive are useful for porting to Europeana. The emphasis will naturally lie on the descriptive metadata. Also part of the provenance data and the rights data should be made available to Europeana, to ensure proper acknowledgment and use of the archive data within Europeana.

⁸⁰ <http://www.loc.gov/standards/premis/guidelines-premismets.pdf> (checked 24 June 2010)

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

With respect to the metadata interoperability with Europeana we conclude the following:

- Although the EDM (see page 15) was not finalized at the time of writing, it appears to be a big step forward. The fact that the original metadata are not “lost in translation” is an essential feature: a *conditio sine qua non*.
- The fact that the notion of aggregations in the EDM allows for part-whole relations between Europeana objects fulfils an impotent requirement of the a/v archives. It allows annotation of video fragments, for example. Also, the clear separation of the work/object/program from (possibly multiple) digital representations is a crucial feature.

EDM is not by itself sufficient for interoperability between Europeana and a/v archive metadata. Support is needed for at least two other aspects. Firstly, as EDM is based on a Web-based formalism, there needs to be a clear scheme for generating and using URIs for a/v fragments. The work of the W3C media fragment working group (see page 23) provides a solution for this and we recommend that PrestoPRIME uses this approach generating interoperable fragment identifiers. Secondly, a/v archives typically do not use Dublin Core as their metadata format. For interoperability with Europeana the a/v archive metadata need to be defined as specialisations of the DCMI terms. The work of the W3C Media Annotation Working Group (MAWG) offers considerable help here (see page 20) and we therefore recommend that PrestoPRIME actively deploys this approach. The annotation ontology specified by MAWG provides a clear route for the specialisations needed to make TV-archive metadata available to Europeana. This includes formats like EBU Core (page 24) and others.

Summarising, the metadata standards described in this deliverable provide the base-level information needed for interoperability with Europeana, assuming only a limited set of metadata is involved in this process.

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

Glossary

Term	Definition
API	Application Programming Interface
BIM	binary format for MPEG-7 data
CDP	Core Description Profile in MPEG-7
CDWA	Categories for the Description of Works of Art
CIDOC CRM	Comité international pour la documentation - Conceptual Reference Model
DAVP	Detailed Audiovisual Profile in MPEG-7
DCMI	Dublin Core Metadata Initiative
DDL	Description Definition Language
DNX	DPS Normalized XML (DPS stands for – Digital Preservation System)
DOI	Document Type Definition
DROID	Digital Record Object Identification
DTD	Document Type Definition
DC	Dublin Core
EAD	Encoded Archival Description
EBU	European Broadcast Union
EDM	Europeana Data Model
EDOB	Editorial Object Documents
ESE	Europeana Semantic Elements
GDFR	Global Digital Format Registry
GTAA	Gemeenschappelijke Thesaurus Audiovisuele Archieven
ISO	International Organization for Standardization
KLV	Key-Length-Value
LIDO	Lightweight Information Describing Objects
LOC	Library of Congress
MAWG	Media Annotation Working Group
MDS	Multimedia Description Schemes in MPEG-7
METS	Metadata Encoding & Transmission Standard
MPEG-7	Multimedia Content Description Interface
MVCO	Media Value Chain Ontology for MPEG-21
MXF	Material Exchange Format

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

Metadata Interoperability	Exchange of metadata where the interpretation of the data remains unchanged
NDHA	National Digital Heritage Archive (of New Zealand)
NLNZ	National Library of New Zealand
OAI-ORE	Open Archives Initiative – Object Reuse and Exchange
PREMIS	Preservation Metadata Implementation Strategies
RDF	Society of Motion Picture and Television Engineers
SGML	Society of Motion Picture and Television Engineers
SKOS	Society of Motion Picture and Television Engineers
SMP	Simple Metadata Profile in MPEG-7
SMPTE	Society of Motion Picture and Television Engineers
TeM	textual format for MPEG-7 data
TVA	TV-Anytime
UDFR	Unified Digital Formats Registry
UDP	User Description Profile in MPEG-7
URI	URI
VRA	Visual Resources Association
XML	Extensible Markup Language
W3C	World Wide Web Consortium

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

References

- [BS06] Werner Bailer and Peter Schallauer. The Detailed Audiovisual Profile: Enabling Interoperability between MPEG-7 based Systems. In 12th International MultiMedia Modelling Conference (MMM'06), pages 217–224, Beijing, China, 2006.
- [BB09] M. Bober and P. Brasnett. MPEG-7 visual signature tools. In Proc. IEEE International Conference on Multimedia and Expo, Jun. 2009.
- [CDP] Information technology – Multimedia content description interface – Part 9: Profiles and levels. ISO/IEC 15938-9:2005, 2005.
- [CUNDIFF] Morgan Cundiff, Presentation: An Introduction to METS; Network Development and MARC Standards Office at the Library of Congress;
http://www.loc.gov/standards/mets/presentations/cat_dir/cat_mets.ppt
- [DMS-1] Material Exchange Format (MXF) -- Descriptive Metadata Scheme-1, SMPTE 380M, 2004.
- [DEBOLE2009] Franca Debole et al.; European Film Gateway public deliverable D2.2 Common interoperability schema for archival resources and filmographic descriptions, report on the common interoperability schema; June 2009;
http://www.europeanfilmgateway.eu/downloads/D2%20Common_Interoperability_Schema_final.pdf (checked 11 June 2010)
- [EBU2009] EBU-TECH 3293: EBU Core Metadata Set (EBU Core); Geneva; July 2009;
http://tech.ebu.ch/docs/tech/tech3293v1_1.pdf (checked: 17 December 2009)
- [TVA1] ETSI TS 102 822-3-1 V1.3.1 Technical Specification. Broadcast and online services: Search, select, and rightful use of content on personal storage systems ("tv anytime"); part 3: Metadata; subpart 1: Phase 1 metadata schemas, June 2005.
- [TVA2] ETSI TS 102 822-2 V1.3.1. Broadcast and online services: Search, select, and rightful use of content on personal storage systems ("tv anytime"); part 2: System description, June 2005.
- [FMS07] James Frew, Dominic Metzger, Peter Slaughter. "Automatic capture and reconstruction of computational provenance." *Concurrency and Computation: Practice and Experience*, 20(5):485-496, Aug. 2007.
- [GM09] Paul Groth, Luc Moreau. "Recording Process Documentation for Provenance," *IEEE Transactions on Parallel and Distributed Systems*, pp. 1246-1259, September, 2009.
- [Hab07] Tom Habing. "Integrating PREMIS and METS", *PREMIS Tutorial, Implementers' Panel*, Washington, DC, Jun. 2007. URL:
<http://www.loc.gov/standards/premis/PremisPanel-habing.ppt>
- [Hoe09] Nancy J. Hoebelheinrich. "PREMIS & Geospatial Resources", *PREMIS Implementation Fair*, San Francisco, CA, Oct. 2009. URL:
http://www.loc.gov/standards/premis/pif-presentations/PREMIS_ImplementFairnjh.ppt
- [HUNTER] J. Hunter, "Working Towards MetaUtopia - A Survey of Current Metadata Research", *Library Trends*, Organizing the Internet, Edited by Andrew Torok, 52(2), Fall 2003
- [Mor07] Moreau, L., Freire, J., Futrelle, J., McGrath, R., Myers, J. and Paulson, P. *The Open Provenance Model*. Technical Report, ECS, University of Southampton, 2007.

D2.2.2 Metadata Models, Interoperability Gaps and Extensions to Preservation Metadata Standards

[MPEG7] MPEG-7, Multimedia Content Description Interface, ISO/IEC 15938, 2001.

[MPF] Metadata production framework specifications (v. 2.0.2E). Technical report, NHK Science and Technical Research Laboratories, Dec. 2008.

<http://www.nhk.or.jp/strl/mpf/english/index.htm>.

[MVCO] Information technology -- Multimedia framework (MPEG-21) -- Part 19: Media Value Chain Ontology. ISO/IEC 21000-19:2010

[MXF] Material Exchange Format (MXF) – File Format Specification (Standard), SMPTE 377M, 2004.

[MXFXML] Material Exchange Format — XML Encoding for Metadata and File Structure Information, SMPTE 434-2006.

[NLNZ2003] Metadata Standards Framework – Preservation Metadata (Revised); June 2003; National Library of New Zealand;

<http://www.natlib.govt.nz/downloads/metaschema-revised.pdf> (checked 10 December 2009)

[PREMIS2008] PREMIS Editorial Committee chaired by Rebecca Guenther; PREMIS Data Dictionary for Preservation Metadata; version 2.0; March 2008;

<http://www.loc.gov/standards/premis/v2/premis-2-0.pdf>

[PREMISMETS2008] Guidelines for using PREMIS with METS for exchange, Revised September 17, 2008; <http://www.loc.gov/standards/premis/guidelines-premismets.pdf>; (checked: 10 December 2009)

[RP210] Metadata Dictionary Registry of Metadata Element Descriptions. SMPTE RP210.11, 2008.

[THOMPSON2003] Dave Thompson, Sam Searle; Preservation Metadata: Pragmatic First Steps at the National Library of New Zealand; D-Lib Magazine April 2003; Volume 9 Number 4; ISSN 1082-9873;

<http://www.dlib.org/dlib/april03/thompson/04thompson.html> (checked: 15 December 2009)

[Woo07] Deborah Woodyard-Robinson. "Implementing the PREMIS data dictionary: a survey of approaches." Report for The PREMIS Maintenance Activity, Jun. 2007. URL: <http://www.loc.gov/standards/premis/implementation-report-woodyard.pdf>

