

Introduction to Digital Libraries

hussein suleman
uct cs honours 2004

OAI Protocol for Metadata Harvesting

What is the OAI?

- What is the Open Archives Initiative (OAI)?
 - Organisation dedicated to solving problems of digital library interoperability by defining simple protocols, most recently for the exchange of metadata.
- What is the Protocol for Metadata Harvesting?
 - Protocol to transfer metadata from a source archive to a destination archive.

Motivation

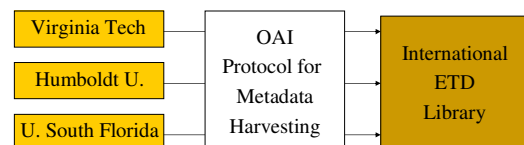
- Existence of some established but independent archives.
- Need for cross-archive services (like search engines).
- Lack of low-cost interoperability technology.
- Experience from past projects such as Dienst.

History

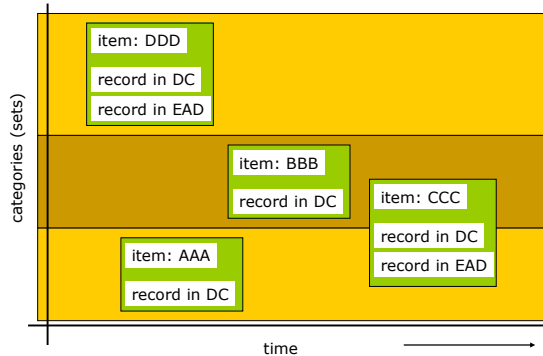
- Santa Fe Meeting – October 1999
 - Santa Fe Convention, January 2000
- Workshops (ACM-DL 2000, ECDL 2000)
- Structuring of the OAI
 - Steering Committee
 - Technical Committee
- Open Days – US/Europe
 - Protocol for Metadata Harvesting v1.0, January 2001
- Minor Update: v1.1 – July 2001
- Version 2.0 – June 2002

Case Study: NDLTD

- Networked Digital Library of Theses and Dissertations
- Made up of multiple independent university-based collections of electronic documents.



Multi-dimensional Data Model



Definitions / Concepts

- Basic Principles
 - What is an Open Archive?
 - Harvesting vs. Federation
 - Metadata vs. Data
 - Data and Service Providers
- Underlying Technology
 - HTTP and XML
 - XML, XML Namespaces and Schema
- Protocol Policies
 - Uniqueness and Persistence
 - What is a record?
 - Multiplicity of Metadata
 - Sets
 - Datestamp, Harvesting and Flow Control

What is an Open Archive ?

- Any WWW-based system that can be accessed through the well-defined interface of the Open Archives Protocol for Metadata Harvesting.
- ... a.k.a. OAI-Compliant Repository
- No implications for:
 - Physical storage of data
 - Cost of data
 - Metadata and data formats
 - Access control to server

Harvesting vs. Federation

- Competing approaches to interoperability
 - Federation is when services are run remotely on remote data (e.g. Federated searching)
 - Harvesting is when data/metadata is transferred from the remote source to the destination where the services are located (e.g. Union catalogues).
- Federation requires more effort at each remote source but is easier for the local system and vice versa for harvesting.
- OAI currently focuses on harvesting.

Metadata vs. Data

- Data refers to digital objects or digital representations of objects.
- Metadata is information about the objects (e.g. title, author, etc.).
- OAI focuses on metadata, with the implicit understanding that metadata usually contains useful links to the source digital objects.

Data and Service Providers

- Data Providers refer to entities who possess data/metadata and are willing to share this with others (internally or externally) via well-defined OAI protocols (e.g., database servers).
- Service Providers are entities who harvest data from Data Providers in order to provide higher-level services to users (e.g. search engines).
- OAI uses these denotations for its client/server model (data=server, service=client).

HTTP and XML

- Protocol for Metadata Harvesting is an almost stateless request/response protocol.
- Requests and responses are sent via the HTTP protocol.
- Requests are encoded as GET/POST operations.
- Responses are well-formed XML documents.

XML Namespaces and Schema

- Consistency and data quality is ensured by using XML Schema descriptions for each possible response.
- XML Namespaces are used where necessary to clearly define which parts of the responses are actual metadata and which support the Protocol for Metadata Harvesting.

Uniqueness and Persistence

- Each record must be uniquely addressable by a distinct identifier.
- Identifiers must be valid URIs
- Example:
 - oai:<archiveId>:<recordId>
 - oai:etd.vt.edu:etd-1234567890
- Each identifier must resolve to a single record and always to the same record (for a given metadata format).

What is a record ?

- A record refers to an independent XML structure that may be associated with digital or physical objects.
- Records are usually associated with metadata, not data.
- OAI advocates harvesting of records, which contain metadata and additional fields to support the harvesting operation.

Sample OAI Record

(note: schema and namespaces have been left out for clarity)

```
<record>
  <header>
    <identifier>oai:jcdl2002.org:tut1</identifier>
    <timestamp>2002-02-03</timestamp>
    <setSpec>tut</setSpec>
  </header>
  <metadata>
    <dc>
      <title>Oldie-but-goodie example</title>
      <creator>Hussein Suleman</creator>
      <language>English</language>
    </dc>
  </metadata>
  <about>
    <metadataID>oai:jcdl2002.org:tut1md</metadataID>
  </about>
</record>
```

Multiplicity of Metadata

- Multiple formats of metadata allowed.
- Dublin Core is mandatory.
- Any other format allowed as long as it has an XML encoding.
- E.g. MARC (Libraries), IMS (Education), ETDS (Theses/Dissertations), RFC1807 (Bibliographies)

Sets

- Protocol mechanism to allow for harvesting of sub-collections.
- No well-defined semantics – depends completely on local data providers.
- May be defined by arrangement between data providers and service providers.
- E.g. Subject areas, years, author names, search queries

Datestamps & Harvesting

- Each record needs a datestamp that indicates its date of creation/modification/deletion.
- Different from dates within the metadata – this date is used only for harvesting
- Can be either YYYY-MM-DD or YYYY-MM-DDThh:mm:ssZ (must be GMT timezone)
- Dates are used to allow for harvesting by date range, thus allowing incremental and continuous transfer of metadata from a data provider to a service provider.

Flow Control

- HTTP “retry-after” mechanism can be leveraged to support server-side delaying of a client’s request.
- Resumption Tokens can be used to return partial results – the client is issued with a token which may be presented to the server to receive more results.

Deletions

- Archives may keep track of deleted records, by identifier and datestamp.
- All protocol result sets can indicate deleted records.
- If deletions are being tracked, this information must be stored indefinitely so as to correctly propagate to service providers with varying harvesting schedules.

Protocol Specifics

- Service Requests
 - Identify
 - ListMetadataFormats
 - ListSets
 - GetRecord
 - ListIdentifiers
 - ListRecords
- Metadata Multiplicity
- Date Ranges
- Resumption Tokens
- Error and Exceptions

Identify

- Purpose
 - Return general information about the archive and its policies
- Parameters
 - None
- Sample URL
 - <http://www.anarchive.org/cgi-bin/OAI?verb=Identify>

Identify - Response

```
Address http://rocky.dlib.vt.edu/~jcdpix/cgi-bin/OAI2.0/beta2/jcd/oai.pl?verb=identify

<?xml version="1.0" encoding="UTF-8" ?>
- <OAI-PMH xmlns="http://www.openarchives.org/OAI/2.0/"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/
http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd"
  <responseDate>2002-05-26T19:28:31Z</responseDate>
  <request verb="Identify">http://rocky.dlib.vt.edu/~jcdpix/cgi-
bin/OAI2.0/beta2/jcd/oai.pl</request>
- <Identify>
  <repositoryName>JCDL 2001 Picture Archive</repositoryName>
  <baseURL>http://rocky.dlib.vt.edu/~jcdpix/cgi-
bin/OAI2.0/beta2/jcd/oai.pl</baseURL>
  <protocolVersion>2.0b2</protocolVersion>
  <adminEmail>jcdpix@rocky.dlib.vt.edu</adminEmail>
  <earliestDatestamp>1970-01-01T00:00:00Z</earliestDatestamp>
  <deletedRecord>no</deletedRecord>
  <granularity>YYYY-MM-DD</granularity>
  <description>
  </Identify>
</OAI-PMH>
```

ListMetadataFormats

- Purpose
 - List metadata formats supported by the archive as well as their schema locations and namespaces
- Parameters
 - identifier – for a specific record (O)
- Sample URL
 - <http://www.anarchive.org/cgi-bin/OAI?verb=ListMetadataFormats>

ListMetadataFormats - Response

```
Address http://rocky.dlib.vt.edu/~jcdpix/cgi-bin/OAI2.0/beta2/jcd/oai.pl?verb=ListMetadataFormats

<?xml version="1.0" encoding="UTF-8" ?>
- <OAI-PMH xmlns="http://www.openarchives.org/OAI/2.0/"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/
http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd"
  <responseDate>2002-05-26T19:27:52Z</responseDate>
  <request verb="ListMetadataFormats">http://rocky.dlib.vt.edu/~jcdpix/cgi-
bin/OAI2.0/beta2/jcd/oai.pl</request>
- <ListMetadataFormats>
  + <metadataFormat>
  - <metadataFormat>
    <metadataPrefix>oai_dc</metadataPrefix>
    <schema>http://www.openarchives.org/OAI/2.0/oai_dc.xsd</schema>
    <metadataNamespace>http://www.openarchives.org/OAI/2.0/oai_dc/</metadataNamespace>
  </metadataFormat>
  </ListMetadataFormats>
</OAI-PMH>
```

ListSets

- Purpose
 - Provide a hierarchical listing of sets in which records may be organised
- Parameters
 - None
- Sample URL
 - <http://www.anarchive.org/cgi-bin/OAI?verb=ListSets>

ListSets – Response

```
Address http://rocky.dlib.vt.edu/~jcdpix/cgi-bin/OAI2.0/beta2/jcd/oai.pl?verb=ListSets

<?xml version="1.0" encoding="UTF-8" ?>
- <OAI-PMH xmlns="http://www.openarchives.org/OAI/2.0/"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/
http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd"
  <responseDate>2002-05-26T19:29:32Z</responseDate>
  <request verb="ListSets">http://rocky.dlib.vt.edu/~jcdpix/cgi-
bin/OAI2.0/beta2/jcd/oai.pl</request>
- <ListSets>
  - <set>
    <setSpec>200105dle</setSpec>
    <setName>JCDL Day Four</setName>
  </setDescription>
  - <oaiddc:dc
    xmlns:oaiddc="http://www.openarchives.org/OAI/2.0/oai_dc/"
    xmlns="http://www.w3.org/2001/XMLSchema-instance"
    xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/oai_dc/
http://www.openarchives.org/OAI/2.0/oai_dc.xsd"
    <description>Pictures taken during JCDL Day Four</description>
  </oaiddc:dc>
  </setDescription>
  </sets>
+ <set>
```

GetRecord

- Purpose
 - Returns the metadata for a single identifier in the form of an OAI record
- Parameters
 - identifier – unique id for record (R)
 - metadataPrefix – metadata format (R)
- Sample URL
 - http://www.anarchive.org/cgi-bin/OAI?verb=GetRecord&identifier=oai:test:123&metadataPrefix=oai_dc

GetRecord - Response

```
Address http://rocky.dlib.vt.edu/~jcdpix/cgi-bin/OAI2.0/beta2/jcd/oai.pl?verb=GetRecord&metadataPrefix=oai_dc
<responseDate>2002-05-26T19:32:54Z</responseDate>
<request verb="GetRecord" metadataPrefix="oai_dc"
  identifier="oai:JCDLPICS:200105dle1">http://rocky.dlib.vt.edu/~jcdpix/cgi-
  bin/OAI2.0/beta2/jcd/oai.pl</request>
- <GetRecord>
- <record>
- <header>
  <identifier>oai:JCDLPICS:200105dle1</identifier>
  <datestamp>2001-06-27</datestamp>
  <setSpec>200105dle</setSpec>
</header>
- <metadata>
- <oaiddc xmlns="http://purl.org/dc/elements/1.1/"
  xmlns:oidc="http://www.openarchives.org/OAI/2.0/oai_dc/"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/oai_dc/
  http://www.openarchives.org/OAI/2.0/oai_dc.xsd">
  <title>01dle1</title>
  <creator>Hussein Suleman</creator>
  <subject>JCDL Day Four</subject>
  <description>Jim French and Unmil Karadkar over lunch</description>
  <publisher>JCDL</publisher>
  <date>2001-06-27</date>
  <type>image</type>
```

ListIdentifiers

- Purpose
 - List headers for all records corresponding to the specified parameters
- Parameters
 - from – start date (O)
 - until – end date (O)
 - set – set to harvest from (O)
 - metadataPrefix – metadata format to list identifiers for (R)
 - resumptionToken – flow control mechanism (X)
- Sample URL
 - http://www.anarchive.org/cgi-bin/OAI?verb=ListIdentifiers&metadataPrefix=oai_dc

ListIdentifiers - Response

```
Address http://rocky.dlib.vt.edu/~jcdpix/cgi-bin/OAI2.0/beta2/jcd/oai.pl?verb=ListIdentifiers&im
<?xml version="1.0" encoding="UTF-8" ?>
- <OAI-PMH xmlns="http://www.openarchives.org/OAI/2.0/"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/
  http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd">
  <responseDate>2002-05-26T19:36:13Z</responseDate>
  <request verb="ListIdentifiers"
    metadataPrefix="oai_dc">http://rocky.dlib.vt.edu/~jcdpix/cgi-
    bin/OAI2.0/beta2/jcd/oai.pl</request>
- <ListIdentifiers>
- <header>
  <identifier>oai:JCDLPICS:200105dle1</identifier>
  <datestamp>2001-06-27</datestamp>
  <setSpec>200105dle</setSpec>
</header>
- <header>
  <identifier>oai:JCDLPICS:200105dle2</identifier>
  <datestamp>2001-06-27</datestamp>
  <setSpec>200105dle</setSpec>
```

ListRecords

- Purpose
 - Retrieves metadata for multiple records
- Parameters
 - from – start date (O)
 - until – end date (O)
 - set – set to harvest from (O)
 - resumptionToken – flow control mechanism (X)
 - metadataPrefix – metadata format (R)
- Sample URL
 - http://www.anarchive.org/cgi-bin/OAI?verb=ListRecord&metadataPrefix=oai_dc&from=2001-01-01

ListRecords - Response

```
Address http://rocky.dlib.vt.edu/~jcdpix/cgi-bin/OAI2.0/beta2/jcd/oai.pl?verb=ListRecords&metadataPrefix=oai_dc
xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/
http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd"
<responseDate>2002-05-26T19:37:37Z</responseDate>
<request verb="ListRecords"
  metadataPrefix="oai_dc">http://rocky.dlib.vt.edu/~jcdpix/cgi-
  bin/OAI2.0/beta2/jcd/oai.pl</request>
- <ListRecords>
- <record>
- <header>
  <identifier>oai:JCDLPICS:200105dle1</identifier>
  <datestamp>2001-06-27</datestamp>
  <setSpec>200105dle</setSpec>
</header>
- <metadata>
- <oaiddc xmlns="http://purl.org/dc/elements/1.1/"
  xmlns:oidc="http://www.openarchives.org/OAI/2.0/oai_dc/"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-
  instance"
  xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/oai_dc/
  http://www.openarchives.org/OAI/2.0/oai_dc.xsd">
  <title>01dle</title>
  <creator>Hussein Suleman</creator>
  <subject>JCDL Day Four</subject>
  <description>Jim French and Unmil Karadkar over
  lunch</description>
```

Metadata Multiplicity

```
- <record>
- <header>
  <identifier>oai:VTETD:etd-3123162539751141</identifier>
  <datestamp>1997-04-22</datestamp>
</header>
- <metadata>
- <rfc1807 xmlns="http://info.internet.isi.edu:80/in-
  notes/rfc/files/rfc1807.txt"
  xsi:schemaLocation="http://info.internet.isi.edu:80/in-
  notes/rfc/files/rfc1807.txt
  http://www.openarchives.org/OAI/rfc1807.xsd">
  <bib-version>1</bib-version>
  <id>etd-3123162539751141</id>
  <entry>1997-04-22</entry>
  <organization>Virginia Polytechnic Institute and State
  University</organization>
  <title>SMA-Induced Deformations In general Unsymmetric
  Laminates</title>
  <tvna>Thesis/Dissertation</tvna>
```

Date Ranges

Address Go

```
<?xml version="1.0" encoding="UTF-8" ?>
- <OAI-PMH xmlns="http://www.openarchives.org/OAI/2.0/"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/
  http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd">
  <responseDate>2002-05-26T19:41:16Z</responseDate>
  <request verb="ListIdentifiers" metadataPrefix="oa_dc" from="2001-06-
  26" until="2001-06-26">http://rocky.dlib.vt.edu/~jcdlpix/cgi-
  bin/OAI2.0/beta2/jcdl/oaipdf</request>
  - <ListIdentifiers>
  - <header>
  <identifier>oai:JCDLPICS:200102d1b1</identifier>
  <timestamp>2001-06-26</timestamp>
  <setSpec>200102d1b1</setSpec>
  </header>
  - <header>
  <identifier>oai:JCDLPICS:200102d1b2</identifier>
  <timestamp>2001-06-26</timestamp>
  <setSpec>200102d1b2</setSpec>
```

Resumption Token

Address Go

```
<identifier>oai:JCDLPICS:200101d1a9</identifier>
<timestamp>2001-06-26</timestamp>
<setSpec>200101d1a9</setSpec>
</header>
- <header>
  <identifier>oai:JCDLPICS:200101d1a10</identifier>
  <timestamp>2001-06-26</timestamp>
  <setSpec>200101d1a10</setSpec>
</header>
<resumptionToken cursor="0" completeListSize="35">!2001-06-26!
2001-06-26!oai_dc!30</resumptionToken>
</ListIdentifiers>
</OAI-PMH>
```

Errors and Exceptions

Address Go

```
<?xml version="1.0" encoding="UTF-8" ?>
- <OAI-PMH xmlns="http://www.openarchives.org/OAI/2.0/"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.openarchives.org/OAI/2.0/
  http://www.openarchives.org/OAI/2.0/OAI-PMH.xsd">
  <responseDate>2002-05-26T19:43:59Z</responseDate>
  <request verb="ListIdentifiers" metadataPrefix="oa_dc" from="2001-06-
  28" until="2001-06-28">http://rocky.dlib.vt.edu/~jcdlpix/cgi-
  bin/OAI2.0/beta2/jcdl/oaipdf</request>
  <error code="noRecordsMatch">The combination of the values of
  arguments results in an empty set</error>
</OAI-PMH>
```

Implementation Details

- Basic requirements
- Basic program layout
- Object-oriented approaches
- Extensible metadata generation
- Data cleaning
- Caching of results
- Error handling
- Denial-of-service prevention
- Creating resumption tokens

Basic Requirements

- You need a WWW Server ☺
- Protocol may be implemented in many forms.
 - CGI Script (Perl, C++, Java)
 - Java Servlet
 - PHP
- Metadata (e.g. database) access mechanism required.
- See www.openarchives.org for list of publicly available software templates.

Basic Program Layout

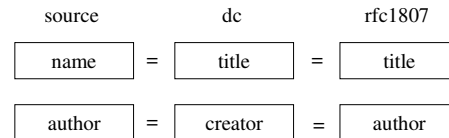
```
parse WWW request to extract parameters
if (verb='Identify')
  ProcessIdentify;
else if (verb='ListMetadataFormats')
  ProcessListMetadataFormats;
else if (verb='ListSets')
  ProcessListSets;
else if (verb='GetRecord')
  ProcessGetRecord;
else if (verb='ListIdentifiers')
  ProcessListIdentifiers;
else if (verb='ListRecords')
  ProcessListRecords;
else
  ReportError ('badVerb');
```

Object-Oriented Approaches

- Cleaner separation of protocol, database access and metadata generation.
- Example approaches
 - Each service request is handled by a object
 - Simpler incremental development
 - Protocol, Database and Metadata are objects
 - Greater portability of code
 - Inheritance from a basic OAI data provider

Metadata Generation

- Approaches
 - Map from source to each metadata format
 - Use crosswalks (maybe XSLT) to generate additional formats.



Data Cleaning

- Escape special XML characters.
- Convert to UTF-8 version of Unicode.
- Convert entity references.
- Remove extraneous whitespace.
- Convert CR/LF for paragraphs.
- URLs
 - `/?#=&.;+` must be encoded as escape sequences

Result Caching

- For multiple requests from many clients or to handle partial result sets.
- Keep temporary tables/files.
- Expire temporary data when no longer needed.
- Is this necessary to handle date-range requests where new items are added to the result set while harvesting is in progress?

Error Handling

- All protocol errors are in XML format
 - badVerb: illegal verb requested
 - badArgument: illegal parameter values or combinations
 - badResumptionToken, cannotDisseminateFormat, idDoesNotExist: parameters are in right format but are not legal under current conditions
 - noRecordsMatch, noMetadataFormats, noSetHierarchy: empty response exception

Denial-of-Service Prevention

- Return only partial results and issue a resumption token for more.
- Use 503 retry-after HTTP errors to have clients try again after a specified back-off time.
- Use access control lists to limit who may access the archive.
- Invoke an explicit delay before sending back results.

Creating resumptionTokens

- Combine from/until/metadataPrefix/set and a record number indicator with delimiters into a sequential token.
For example:
 - from!until!metadataPrefix!set!recordnumber
 - 2000-01-01!2001-01-01!!All!100
- Use a session manager with automatic expiry.
For example:
 - vtetd14june10amsession12

Tools for Testing

- Repository Explorer
 - Interactive Browsing
 - Testing of parameters
 - Multiple views of data
 - Multilingual support
 - Automatic test suite
- OAI Registry
- XML Schema Validator

RE Interactive Browsing

RE Parameter Testing

RE Browsing

RE Browsing

List of Sets

Click on the link to list the contents

[JCDDL Day Four](#)

set description:

dc:
description: Pictures taken during JCDDL Day Four

[JCDDL Banquet](#)

set description:

dc:
description: Pictures taken during JCDDL Banquet

[JCDDL Day Three](#)

RE Browsing

List of Record Identifiers

Select a link to view more information

header:

identifier : oai:JCDLPICS:200105d1e1
datestamp : 2001-06-27
setSpec : 200105d1e

[\[display record in Dublin Core\]](#) [\[display metadata formats\]](#)

header:

identifier : oai:JCDLPICS:200105d1e2
datestamp : 2001-06-27
setSpec : 200105d1e

[\[display record in Dublin Core\]](#) [\[display metadata formats\]](#)

RE Browsing

List of Metadata Formats

Click on the link to view schema

Prefix=[dc2]

Namespace=[http://www.openarchives.org/OAI/2.0/oai_dc/]

Schema=[http://www.openarchives.org/OAI/2.0/oai_dc.xsd]

[Not a standard OAI metadata name] [\[display record\]](#)

Prefix=[oai_dc]

Namespace=[http://www.openarchives.org/OAI/2.0/oai_dc/]

Schema=[http://www.openarchives.org/OAI/2.0/oai_dc.xsd]

[\[display record\]](#)

RE Browsing

List of Fields

header:

identifier : oai:JCDLPICS:200105d1e1
datestamp : 2001-06-27
setSpec : 200105d1e

metadata:

dc:

title: Oldie1
creator: Hussein Suleman
subject: JCDL Day Four
description: Jim French and Umil Karadkar over lunch
publisher: JCDL
date: 2001-06-27
type: image
format: image/jpeg
identifier: <http://rocky.dlib.vt.edu/~jcdipix/pictures/200105d1e/Oldie1.jpg>
language: en-us
relation: <http://www.jcdl.org>
rights: unrestricted

RE Multiple views of data

Raw XML Output

```
<?xml version="1.0" encoding="UTF-8"?>
<OAI-PMH xmlns="http://www.openarchives.org/OAI/2.0/" xmlns:xsi="
<responseDate>2002-05-26T19:59:35Z</responseDate>
<request verb="GetRecord" metadataPrefix="oai_dc" identifier="oai:
<GetRecord>
<record>
<header>
<identifier>oai:JCDLPICS:200105d1e1</identifier>
<datestamp>2001-06-27</datestamp>
<setSpec>200105d1e</setSpec>
</header>
<metadata>
<oaid:dc xmlns="http://purl.org/dc/elements/1.1/" xmlns:oaid="
  <title>Oldie1</title>
  <creator>Hussein Suleman</creator>
  <subject>JCDL Day Four</subject>
```

RE Multilingual Support



"Explorer" Version - 1.44 - Protokollversion - 1.01.1/2.0b2 : A

Hier koennen Sie interaktiv die OAI-Kompatibilitaet ihrer Archive verifizieren. [\[Bitte hier klicken um D](#)

JavaScript is required

Bemerkung Um HTTP-Fehler zu vermeiden, warten Sie bitte bis eine Seite fertig geladen ist, bevor Sie klicken

Bitte geben Sie die URL zu dem OAI Interface ein (alles vor dem "?") oder waehlen Sie ein vordefiniertes Tabelle

RE Automatic Test Suite



Open Archives Initiative - Repository Explorer

explorer version - 1.44 / protocol version - 2.0b2 : May 2002

Open Archives Initiative :: Protocol for Metadata Harvesting v2.0b2
RE Protocol Tester 1.44 :: Virginia Tech DURL :: May 2002

```
(1) Testing : Identify
URL : http://rocky.dlib.vt.edu/~jcdipix/cgi-bin/OAI2.0/beta2/jcdl/oai.pl?verb=Identify
Test Result : OK
---- [ Repository Name = JCDL 2001 Picture Archive ]
---- [ Protocol Version = 2.0b2 ]
---- [ Base URL = http://rocky.dlib.vt.edu/~jcdipix/cgi-bin/OAI2.0/beta2/jcdl/oai.pl ]
---- [ Admin Email = jcdipix@rocky.dlib.vt.edu ]
---- [ Granularity = YYYY-MM-DD ]

(2) Testing : Identify (illegal parameter)
URL : http://rocky.dlib.vt.edu/~jcdipix/cgi-bin/OAI2.0/beta2/jcdl/oai.pl?verb=Identify&
Test Result : OK

(3) Testing : ListMetadataFormats
URL : http://rocky.dlib.vt.edu/~jcdipix/cgi-bin/OAI2.0/beta2/jcdl/oai.pl?verb=ListMetad
Test Result : OK
---- [ Sample Metadata Format = dc2 ]
```

RE Error in Response

Archive Self-Description

Repository Name	Virginia Tech Electronic Thesis and Dissertation Collection
Base URL	http://oai.dlib.vt.edu/30/-houssein/cgi-bin/NDLTDErr1/VTETD.pl
Protocol Version	1.0
Error Missing field: <Identify/><adminEmail>	
Other Information	<pre> description: oai-identifier: scheme: oai repositoryIdentifier: VTETD delimiter: ; sampleIdentifier: oai:VTETD:etd-17111020975060 description: eprints: content: text: Theses and Dissertations produced by students at Virginia metadataPolicy: text: Metadata may be used by commercial and non-commercial user dataPolicy: text: Full texts are individually tagged and the rights statement </pre>

RE Error in XML

Internet Explorer version - 1.1 ; protocol version - 1.0 ; April 2001

<http://oai.dlib.vt.edu/-houssein/cgi-bin/NDLTDErr1/VTETD.pl?verb=Identify>

XSD Schema/Instance Validation Error !

Errors in XML instance

```

<?xml version="1.0"?>
<xsxy docElct=" (http://www.openarchives.org/OAI/1.0/OAI_Identifier) Identify" instanceAses
<importAttempt URI="http://oai.dlib.vt.edu/OAI/1.0/OAI_Identifier.xsd" namespace="http://
<importAttempt URI="http://oai.dlib.vt.edu/OAI/1.0/OAI_Identifier.xsd" namespace="http://
<importAttempt URI="http://oai.dlib.vt.edu/OAI/1.0/OAI_Identifier.xsd" namespace="http://www.open
<invalid char="4" code="cvc-complex-type.1.2.4" line="11" resource="file:///tmp/111020975060
<fmo
<node id="1">
<edge dest="2" label="(http://www.openarchives.org/OAI/1.0/OAI_Identifier):responseDate"
</node>
<node id="2">
<edge dest="3" label="(http://www.openarchives.org/OAI/1.0/OAI_Identifier):requestURL"/>
</node>
<node id="3">
          
```

OAI Registry



The Open Archives Initiative

Registering as a Data Provider

Data providers who support the OAI protocol may choose to list their repository in the OAI registry. The goals of the registry are:

- Provide a publicly accessible list of OAI conformant repositories, making it easy for service providers to discover repositories from which metadata can be harvested.
- Provide a mechanism for data providers to ensure their conformance with the OAI protocol specification.
- Provide a means for the OAI to monitor use of the protocol and plan future activities and strategies.

This page allows you to register your repository by entering your [BASE-URL](#) in the text box at the bottom of this page. *Before doing that, please read all of this instruction page so you understand what registration means and the choices you have.*

[Consequences of Registration](#)
[Protocol Testing](#)
[Conformance Testing](#)

OAI Registry



The Open Archives Initiative

List of Registered, OAI-Conformant Repositories

This application allows you to browse the current list of OAI conformant repositories. Currently there are 29 such repositories. The table may be sorted either by the [OAI Repository Identifier](#) or by the [Repository Name](#).

You may retrieve information about an OAI repository by selecting one of the rows in the following table. You may view the registration record from the database; alternatively, if your browser can render XML, you may issue the [Identify request](#) to the selected repository and receive the current XML response.

Sort repositories by:

Repository Name

OAI Identifier

view registration record

issue Identify request

OAI Repository Identifier

- celebration
- andc
- arXiv
- CDLCLAS
- ...

Repository Name

- A Celebration of Women Writers
- Alaska Native Language Center
- arXiv
- California International and Area Studies Digital Repository
- ...

Service Providers

- How to Harvest
- Policies
- Intermediate systems
- Case Study: ARC
- Case Study: NDLTD

How To Harvest

- Identify to get basic information.
- ListIdentifiers, followed by ListMetadataFormats for each record and then GetRecord for each id/metadata combination.
 - No. of short HTTP requests = $1+n+n \times m$
 n =no. of identifiers, m =no. of metadata formats
- ListRecords for each metadata format required.
 - No. of long HTTP requests = m
 m =no. of metadata formats

Policies

- ❑ Use schedule for harvesting regularly.
- ❑ Store date when last harvested (before you start).
- ❑ Use a two day overlap (or one day if your archive uses proper UTC timestamps).
 - New items may be added for the current day.
 - Timezones create up to a day of lag if you ignore them.
 - If the source uses correct UTC timestamps and second granularity then only 1 second of overlap is needed!
- ❑ Each time a record is encountered, erase previous instances.

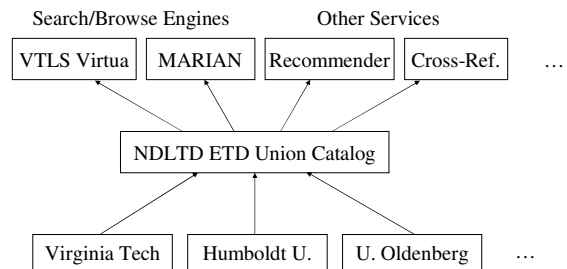
Intermediate Systems

- ❑ Both a data provider and service provider.
- ❑ All harvested data must have the timestamps updated to the date on which the harvesting was done.
- ❑ Identifiers retain their original values.
- ❑ Note: Consistency in the source archive propagates, but so does inconsistency!

8.5. Case Study: ARC



Case Study: NDLTD



Links

- ❑ Open Archives Initiative
 - <http://www.openarchives.org>
- ❑ OAI Metadata Harvesting Protocol
 - <http://www.openarchives.org/OAI/openarchivesprotocol.htm>
- ❑ Virginia Tech DLRL OAI Projects
 - <http://www.dlib.vt.edu/projects/OAI/>
- ❑ Repository Explorer
 - http://purl.org/net/oai_explorer
- ❑ NDLTD
 - <http://www.ndltd.org>

More Links

- ❑ ARC Cross-Archive Search Service
 - <http://arc.cs.odu.edu/>
- ❑ XML Schema Validator
 - <http://www.w3.org/2001/03/webdata/xsv>
- ❑ Dublin Core Metadata Initiative
 - <http://www.dublincore.org>
- ❑ E-Prints DL-in-a-box
 - <http://www.eprints.org>
- ❑ XML Tools at W3C
 - <http://www.w3.org/XML/#software>