open
archives
forum

open
archives forum

# 2nd Technical
# Validation Questionnaire
## - interim results -

Birgit Matthaei

Humboldt-University, Berlin, Germany

Electronic Publishing Group

Computer- and Mediaservice

birgit.matthaei@cms.hu-berlin.de

# Why this technical questionnaire?

➢ **1st Technical Validation Questionnaire**

- provide an overview on status, experiences and future plans belonging OAI implementations of participants of the 1st OAForum Workshop

- target group: workshop participants
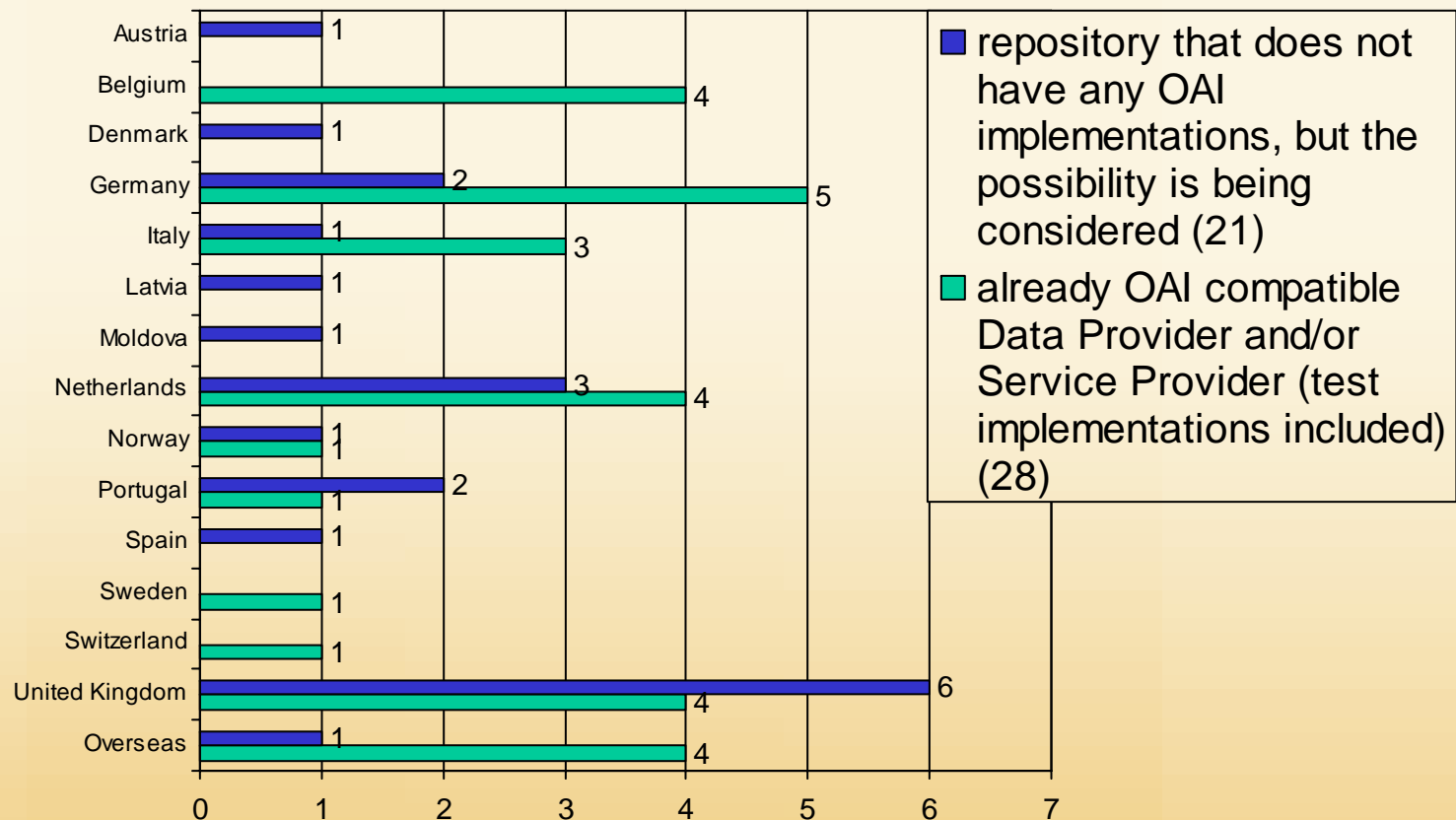
➢ **High Interest, Feedback**

- to collect experiences of a broader spectrum

- to learn more about starting conditions of planned implementations

  ° Is there large common ground?

  ° Are requirements so individual that it will be necessary for many isolated solutions to be developed?

  ° Should tools and protocols correspond more than now to the needs of different communities?

IST- 2001-320015

Birgit Matthaei, 28. March 2003, Berlin, 3rd OAForum Workshop: Networking Multimedia Resources
Humboldt-University, Berlin, Germany - Electronic Publishing Group - CMS / University Library

# What are the Goals?

➢ **Extended 2nd Questionnaire**

- extended questions + target audience + duration
- new subdivision in two questionnaires
  - ° technical presuppositions of those, which have not yet integrated OAI-PMH
  - ° experiences of implementers

➢ **to get information about**

- used software
- implementation costs
- offered spectrum and interoperability
- experiences and expectations
  - → in different communities
  - → in different countries

➢ **to share experiences and information about technical issues related to open archives**
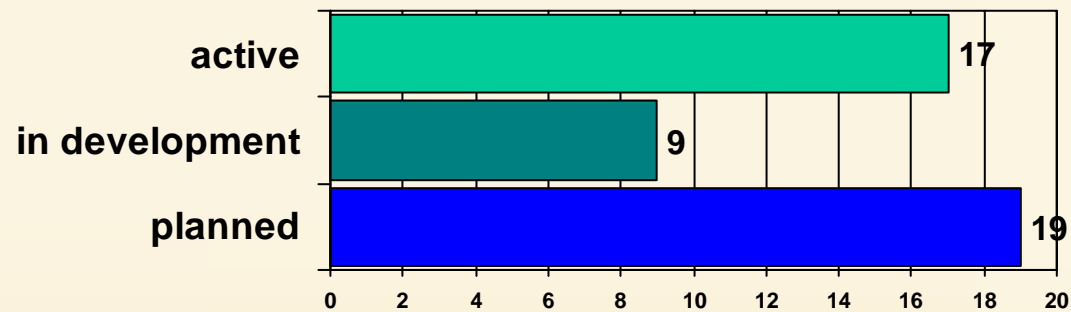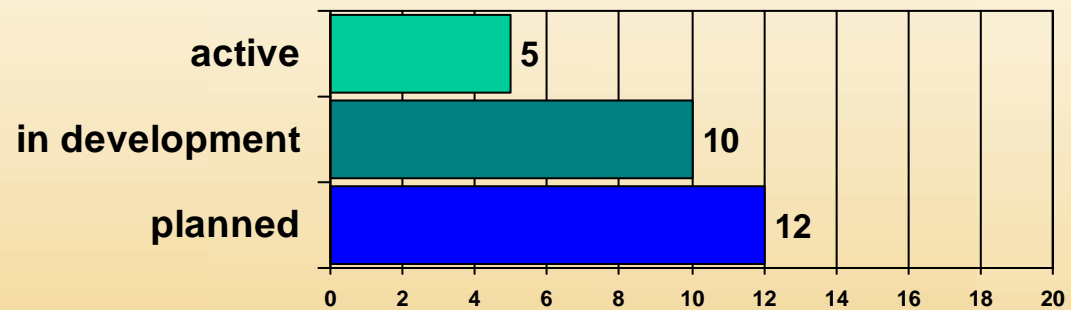
# Who participated to date?

➢ **Countries**

Birgit Matthaei, 28. March 2003, Berlin, 3rd OAForum Workshop: Networking Multimedia Resources
Humboldt-University, Berlin, Germany - Electronic Publishing Group - CMS / University Library

# Who participated to date?

➢ **Data Provider**

| | active | in development | planned |
|---|---|---|---|
| value | 17 | 9 | 19 |

➢ **Service Provider**

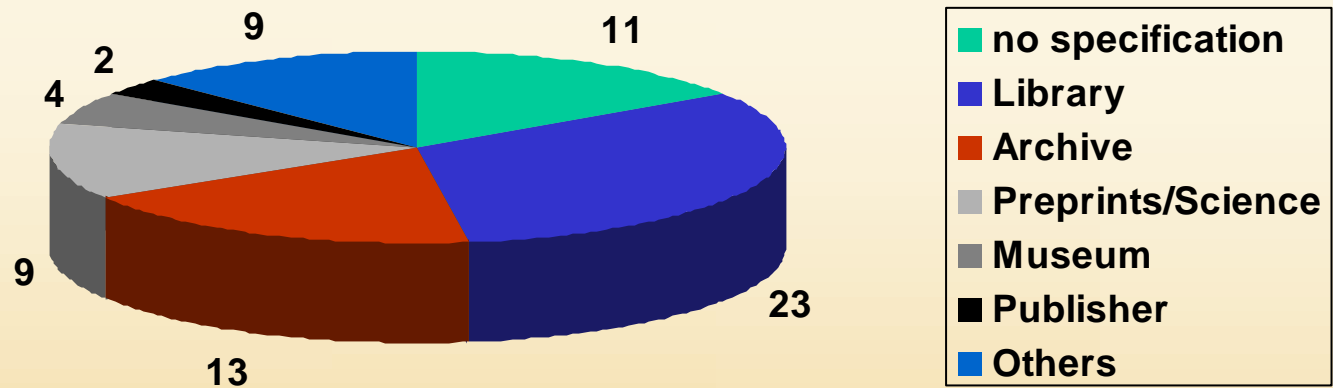| | active | in development | planned |
|---|---|---|---|
| value | 5 | 10 | 12 |

30 % of active DP are also SP

41 % of active DP plan or still develop SP implementations

# Who participated to date?

➢ **Communities**

Multiple answers possible



Legend:
- no specification
- Library
- Archive
- Preprints/Science
- Museum
- Publisher
- Others

Values shown: 11, 23, 13, 9, 4, 2, 9

# Used Software

➤ **Technical infrastructure before OAI-Implementation**

- not many statements
  to Interface and Collection Systems

- dominant programming languages:
  ° Perl, XML, also Java, PHP

- dominant databases:
  ° MySQL, Oracle

➤ **Almost no one changed existing software tools to be OAI compatible**

# Used Software

➢ **Implementations to be OAI compatible**

- about 60 % of the used tools were self-developed by both Data- and Service Providers

    - most of them make their developments and the source code available for others

    - dominant programming languages: Java, Perl, PHP, also XML

- tools like PERL implementations, OAI Cat, EPrints, and OAI Harvester were mentioned 3 or 4 times each

➔ list of OAI-PMH software: http://www.openarchives.org/tools/

➢ **Necessary Know How:
Data- & Service Provider**

focused on various combinations of the following **five competence fields**:

- system administration (UNIX | Linux)

- web server configuration (Apache)

- knowledge on Databases and SQL
  (MySQL | Sybase | Oracle)

- programming skills
  (Perl | Java | PHP | Servlets | CGI | XML)

- experiences with metadata

# Implementation Costs

➢ **Time and Manpower**

- **implementations of OAI-specifications:**
    - ° 75% concluded within a quarter by one programmer (span: from 2 to 750 personal days per month)

- **reasons for few bigger expenditures:**
    - ° context of bigger research projects
    - ° construction of archives
    - ° processing of bigger data amounts

- **further maintenance for a stable protocol:**
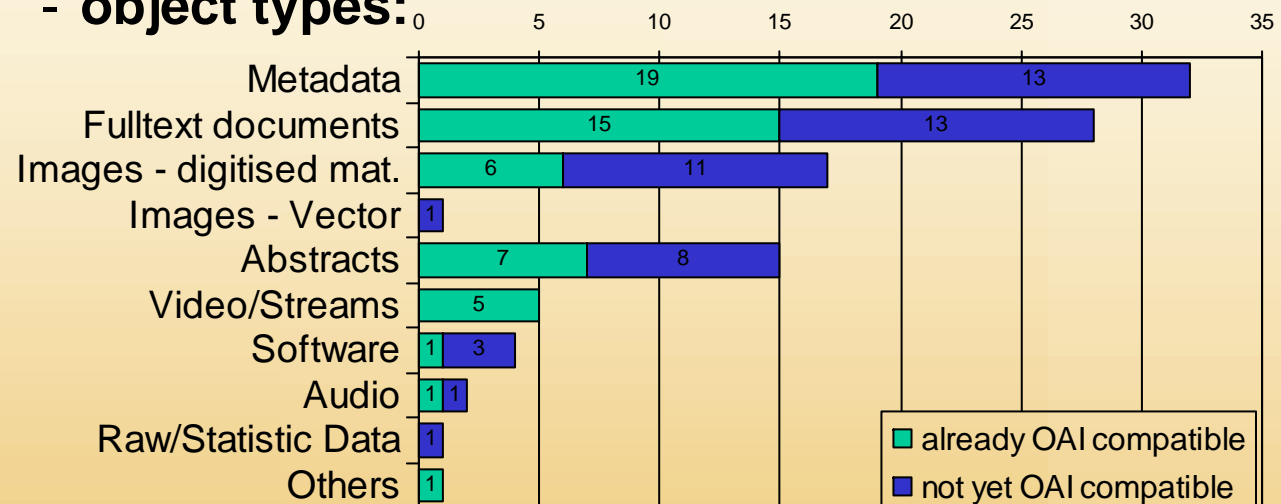    - ° max. 25, mostly 1 personal day per month

# Implementation Costs

➢ **Expectations of those who haven't implemented yet**

- **implementations of OAI-specifications** (same):
  - ° concluded within a quarter by one programmer

- **further maintenance for a stable protocol** (higher):
  - ° up to 40 personal days per month

- No specific trend recognizable with expectations if
  - ° data structures suggested by the OAI-PMH are easy to integrate in existing infrastructure
  - ° the adaption of the data to the OAI-PMH will be expensive
  - ° the preparation of the data for an internet usage will be expensive
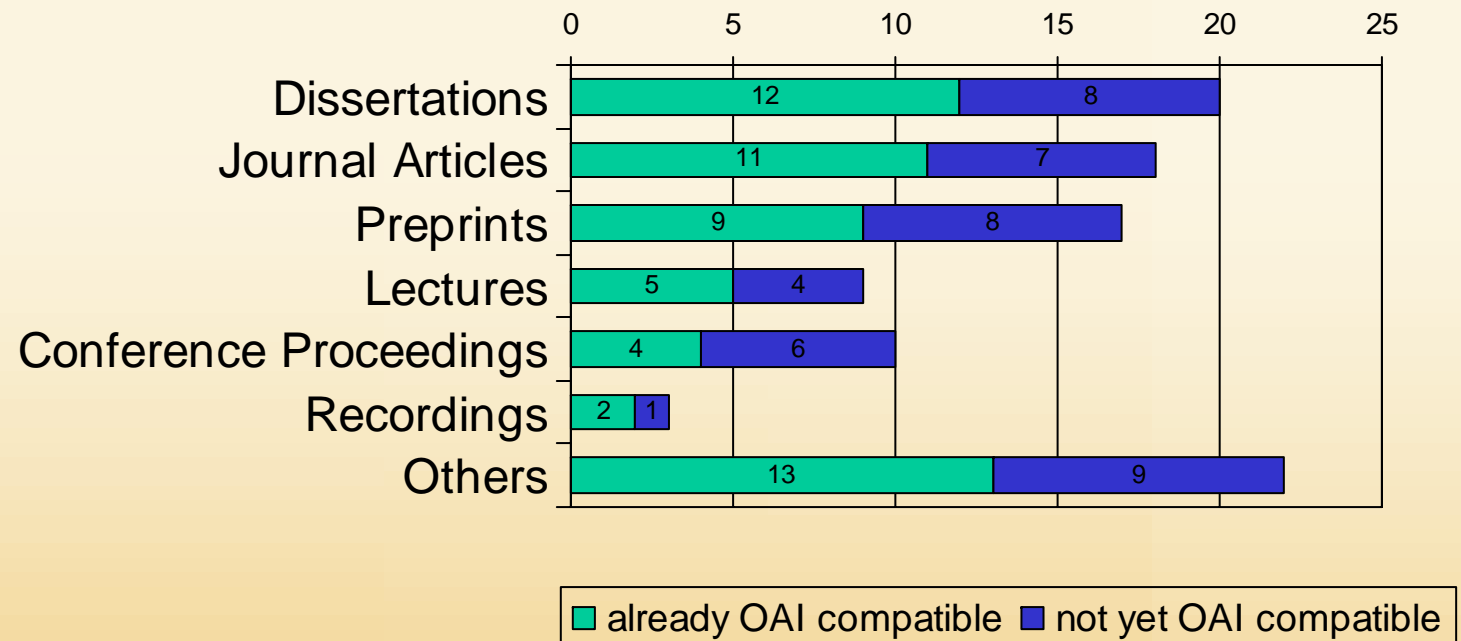
# Offered spectrum - DP

➢ **Offers of Data Providers**

- **number of documents:**
  - ° between 5 and several million documents

- **storage space:**
  - ° between 1 megabytes and 2 Terabyte.

- **object types:**

| Object type | already OAI compatible | not yet OAI compatible |
|---|---|---|
| Metadata | 19 | 13 |
| Fulltext documents | 15 | 13 |
| Images - digitised mat. | 6 | 11 |
| Images - Vector | | 1 |
| Abstracts | 7 | 8 |
| Video/Streams | 5 | |
| Software | 1 | 3 |
| Audio | 1 | 1 |
| Raw/Statistic Data | | 1 |
| Others | 1 | |

Legend: ▇ already OAI compatible ▇ not yet OAI compatible

- **Content types**



Bar chart showing content types, measured on a scale from 0 to 25:

| Content type | already OAI compatible | not yet OAI compatible |
|---|---|---|
| Dissertations | 12 | 8 |
| Journal Articles | 11 | 7 |
| Preprints | 9 | 8 |
| Lectures | 5 | 4 |
| Conference Proceedings | 4 | 6 |
| Recordings | 2 | 1 |
| Others | 13 | 9 |

☐ already OAI compatible ☐ not yet OAI compatible

IST- 2001-320015

Birgit Matthaei, 28. March 2003, Berlin, 3rd OAForum Workshop: Networking Multimedia Resources
Humboldt-University, Berlin, Germany - Electronic Publishing Group - CMS / University Library

# Offered spectrum - DP

- **Metadata formats**

|        | 0 | 5 | 10 | 15 | 20 |
|--------|---|---|----|----|----|

Dublic Core simple — already OAI compatible: **14**, not yet OAI compatible: **5**

Dublin Core qualified — already OAI compatible: **7**, not yet OAI compatible: **4**

MARC 21 — already OAI compatible: **4**, not yet OAI compatible: **2**

UNIMARC — already OAI compatible: **3**, not yet OAI compatible: **2**

MAB — already OAI compatible: **1**, not yet OAI compatible: **2**

EAD — already OAI compatible: **1**, not yet OAI compatible: **2**

Others (single mentioned) — already OAI compatible: **11**, not yet OAI compatible: **2**

☐ already OAI compatible   ☐ not yet OAI compatible

Single mentioned formats:

Dublin Core Library Profile, DiTeD, CEOS CIP, AMF, RIS, MODS, METS, SPECTRUM, TEI, internal format, self developed

# Offered spectrum - DP

- **Dissemination**

  ° more than half of the Data Providers are offering all parts or rather extracts of the documents

  ° if the openness of the OAI interface is reduced due to several reasons, people use two limitation strategies:

    • access control
      (control of the IP-addresses, licensing, agreements)

    • limitation of the data output

# Offered spectrum - SP

➤ **Kind of Services**

- OAI-Service / Portal

- local or community specific services

- searching and browsing for information

- search in different sources through one search interface

- cross-linking, annotations, harvesting

- workspace for managing documents and metadata, collaboration within groups of users

- document management

# Offered spectrum - SP

➢ **Stategies to process with harvested data from DP**

- use no provenance information

- filter harvester output and load local database

- strategies to include information about DP
  in data output:
  - ° when a metadata record is found, the user can also browse information on the archive the record came from
  - ° queries against the portal return data sets as harvested, including information about the original data provider
  - ° provenance information is encoded in the identifier

# Experiences - DP

➢ **Importance / Advantages of OAI**

- provide additional services to existing services
- replace existing services through OAI interface
- better retrieval, make Metadata exchange available
- share scientific knowledge, harvest other knowledge databases, cross-search in institutional assets
- major dissemination of researchers' results
- simple and cheap in implementation
- easy adaption for project internal usage
- simple to implement facility of exchanging metadata in comparison to more complex protocols

➔ „provide access to all of human knowledge"

➔ „nothing other than political expediency"

➢ **Problem: Standardisation**

- heterogeneity of the content of the metadata records requires the service provider to expend a lot of effort in normalizing the data in order to make it more comparable and usable

    ° could be done at lesser cost by the individual data provider

    ° development of middleware tools that service providers could use for data normalization

# Experiences - SP

➢ **Future Plannings**

- extend search & browse functions

- export in other formats such as XML

- document delivery services, print on demand

- collaboration environment for users and groups of users, discussion forums, annotations, awareness

- extend existing services, building distributed services

- establish an exchange of different library catalogues and the integration into a virtual union catalogue for the whole country

- create a single catalogue of all library's catalogues: library opac, archives database, image database, Internet gateways

# Useful information sources

➢ **Problems to find useful informations?**

- Many of those who haven't implemented yet made the experience that it is laborious to find good informations about metadata and especially technical support
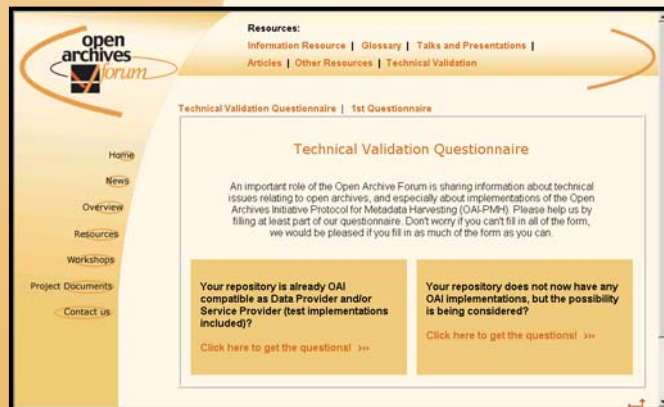- Some asked for a gentle introduction to the protocol („too jargonish")

➢ **Recommendations of the questionnaire participants**

- **Websites**
  - www.openarchives.org
  - www.ndltd.org
  - www.cimi.org
  - www.eprints.org
  - www.rlg.org
  - www.oaforum.org
  - www.ukoln.ac.uk/distributed-systems/jisc-ie/arch/faq/oai
  - http://library.cern.ch/heplw/4/papers/3/
- **Online journals** eg. Ariadne, D-Lib Magazine
  - www.ariadne.ac.uk
  - www.dlib.org
- **Conferences** and workshops
- **Informal discussions** with other gateway managers
- **Test programs** eg. http://oai.dlib.vt.edu/cgi-bin/Explorer/oai2.0/testoai

# Thank You!

➢ **Please contribute!**

- Information about your projects

- Your implementation and usage experience



**Technical Validation Questionnaire**
**http://www.oaforum.org/resources/tecvalq2.php**



**Information Resource Database**
**http://www.oaforum.org/oaf_db/**