

Opening the Archives: OLAC, TRACTOR and the OTA

Martin Wynne
Oxford Text Archive
martin.wynne@ota.ahds.ac.uk

Outline

1. Why did I want to join in with OLAC?
2. What is TRACTOR?
3. What is the Oxford Text Archive?
4. The experiment
5. The results

Martin Wynne, Opening the Archives

2

Standards

- I don't like them;
- But some level of interoperability of metadata is necessary for a centralised resource discovery mechanism;
- Resource encoding standards and metadata standards;
- There are good standards and bad standards.

Martin Wynne, Opening the Archives

3

Bad standards

- Devise elaborate standards in the abstract without reference to practicalities;
- Make complex standards with many obligatory categories;
- Use technologies with which the community is not familiar or which are arcane.

Martin Wynne, Opening the Archives

4

Good standards

- Easy to implement, especially making migration from other formats easy;
- Provide tools and support for migration and implementation;
- Make validation easy;
- Give immediate rewards, e.g. a working search engine;
- Use or relate to existing community (and wider) standards.

Martin Wynne, Opening the Archives

5

Portals

UK:

- Resource Discovery Network
- HUMBUL Humanities Hub
- Arts and Humanities Portal

International:

- Internet search engines
- Subject-based meta-archives (e.g. OLAC)

Martin Wynne, Opening the Archives

6

An alternative: IMDI

- More detailed metadata set
- DC is only the lowest common denominator: need more detailed set for describing language resources properly
- But what categories for what sort of resources?
- IMDI really for linguistic field data, not corpora
- Can now be migrated to OLAC

Martin Wynne, Opening the Archives

7

Conclusions so far

- Need a meta-archive so that resources in different archives can more easily be discovered
- OLAC seems to offer a potentially viable model
- Let's try it out with an archive (or two!)

Martin Wynne, Opening the Archives

8

TRACTOR

- A key part of TELRI II
- www.tractor.de
- Resources in *Bulgarian, Croatian, Czech, Dutch, English, Estonian, French, German, Greek, Hungarian, Italian, Latvian, Lithuanian, Romanian, Russian, Serbian, Slovak, Slovenian, Swedish, Turkish, Ukrainian, Uzbek.*
- Mostly corpora; some lexicons and tools
- c. 100 resources

Martin Wynne, Opening the Archives

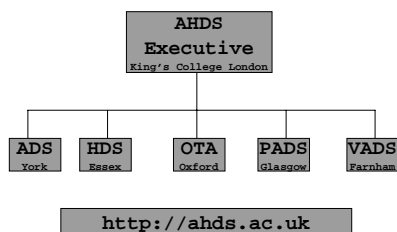
9



Martin Wynne, Opening the Archives

10

Arts and Humanities Data Service structure and services



Martin Wynne, Opening the Archives

11

Oxford Text Archive: Holdings

- 25+ languages, c. 2,500 titles
 - Literary texts
 - Reference works, dictionaries, Bibles
 - Corpora (50+)
- Catalogue (TEI SGML and XML, library cataloguing guidelines followed)
- Full text retrieval in various formats
 - HTML, TEI-Lite, ASCII, RTF, original(?)
- Mirror sites

Martin Wynne, Opening the Archives

12

The experiment: prerequisites

- Structured metadata
- Access to standards and ability to migrate
- Validation tools
- Harvest mechanism

Martin Wynne, Opening the Archives

13

Opening TRACTOR and the OTA to OLAC (1)

Two different archives:

TRACTOR: 99 records, corpora, designed for LE, not much metadata

OTA: 2500 records, mostly literature, some corpora, metadata well-structured but partial

Common points: variable quality of resources and metadata

Martin Wynne, Opening the Archives

14

Process

- Sed, awk, commands and scripts and perl programs to migrate to OLAC-conformant metadata records
- Hack XML file of records
- Validate records
- Send XML file to OLAC

Martin Wynne, Opening the Archives

15

Conversion of TRACTOR metadata (1)

- Parse (reasonably well-structured) HTML file to get metadata categories;
- Re-order (in Excel) and edit to fit OLAC categories (and export to tab-separated text);
- Hack OLAC-conformant XML file from the text.

Martin Wynne, Opening the Archives

16

Conversion of TRACTOR metadata (2)

```
<hr>
<h2>
<a NAME="Bulgarian"></a><i>Bulgarian</i></h2>

<ul>
<p><i><h3>
POS tagged corpus</h3>
2460 Bulgarian sentences marked-up with part of speech
information (BTB-POS Corpus I). The corpus
is in XML format, non-standard with respect to TEI or CES, DTD
is included. Available in three
different encodings of cyrillic letters: ISO 8879:1986, MS
Windows, and Unicode.
<p><i>Resource provider:</i> <a
href="mailto:kivs@bgciet.acad.bg">Kiril Iv. Simov</a>,
Linguistic Modelling Laboratory, Bulgarian
Academy of Sciences, Sofia.
<p><a href="/tractor/resources/SOF2/BTB-POS/">Browse the
files</a>
```

Martin Wynne, Opening the Archives

17

Conversion of TRACTOR metadata (3)

```
SOF201 2001-10-31 "Kiril Iv. Simov, Linguistic Modelling Laboratory, Bulgarian" 2460
sentences unknown 2001 "2460 Bulgarian sentences mar
ked-up with part of speech information (BTB-POS Corpus I). The corpus is in XML format, non-
standard with respect to TEI or CES, DTD is included. Ava
ilable in three different encodings of cyrillic letters: ISO 8879:1986, MS Windows, and Unicode. "
text
ISO 8879:1986 XML with DTD
http://tractor.lham.ac.uk/tractor/resources/SOF2/BTB-POS/ Bulgarian
www.tractor.de TRACTOR User Agreement
BTB-POS Corpus I corpus
SOF202 2001-10-31 "Kiril Iv. Simov, Linguistic Modelling Laboratory, Bulgarian" 2460
sentences unknown 2001 "2460 Bulgarian sentences mar
ked-up with part of speech information (BTB-POS Corpus I). The corpus is in XML format, non-
standard with respect to TEI or CES, DTD is included. Ava
ilable in three different encodings of cyrillic letters: ISO 8879:1986, MS Windows, and Unicode. "
text
ISO 8879:1986 XML with DTD
http://tractor.lham.ac.uk/tractor/resources/SOF2/BTB-POS/ Bulgarian
www.tractor.de TRACTOR User Agreement
BTB-POS Corpus I corpus
SOF203 2001-10-31 "Kiril Iv. Simov, Linguistic Modelling Laboratory, Bulgarian" 2460
sentences unknown 2001 "2460 Bulgarian sentences mar
ked-up with part of speech information (BTB-POS Corpus I). The corpus is in XML format, non-
standard with respect to TEI or CES, DTD is included. Ava
ilable in three different encodings of cyrillic letters: ISO 8879:1986, MS Windows, and Unicode. "
text
ISO 8879:1986 XML with DTD
http://tractor.lham.ac.uk/tractor/resources/SOF2/BTB-POS/ Bulgarian
www.tractor.de TRACTOR User Agreement
BTB-POS Corpus I corpus
```

Martin Wynne, Opening the Archives

18

Conversion of TRACTOR metadata (4)

```

<?xml version="1.0" encoding="UTF-8" ?>
<!DOCTYPE OLAC-Repository SYSTEM "olacrep.dtd">
<OLAC-Repository>
  <identity>
    <repositoryIdentifier>tractorest</repositoryIdentifier>
    <repositoryName>TRACTOR Text Archive</repositoryName>
    <adminEmail>martin.wynne@ota.ahds.ac.uk</adminEmail>
    <collectionDescription>A test set of metadata records from TRACTOR (TELRI Research Archive of Language Resources and
  </collectionDescription>
  </identity>
  <sets>
    <set>
      <setSpec></setSpec>
      <setName></setName>
    </set>
    <set>
      <setSpec></setSpec>
      <setName></setName>
    </set>
    <record>
      <recordId>olactractorSQF2001</recordId>
      <datestamp>2001-10-31</datestamp>
      <header>
        <metadata>
          <dc:lang>"en">
            <contributor>Kiril By Simov, Linguistic Modelling Laboratory, Bulgaria</contributor>
            <coverage>2460 sentences</coverage>
            <creator>unknown</creator>
            <date>2001</date>
            <description>2460 Bulgarian sentences marked-up with part of speech information (RTL-POS Corpus 1). The corpus is in XML
            format, non-standard with respect to TEI or CJK. DTD is included. Available in three different encodings of cyrillic letters: ISO 8879:1986,
            MS Windows, and Unicode. </description>
            <format>text</format>
          </dc:lang>
        </metadata>
      </header>
    </record>
  </sets>
  </OLAC-Repository>

```

Martin Wynne, Opening the Archives

19

Conversion of OTA metadata (1)

```

<!DOCTYPE TEI PUBLIC "-//TEI/DTD TEI Lite 1.0/EN" "tei-lite.dtd">
<TEI>
  <TEIHEADER TYPE="ISBINDER">
    <FILEDESC>
      <TITLESTMT>
        <TITLE TYPE="main">Erec (Electronic resource)</TITLE>
        <AUTHOR>Chok education, de Troyes, 12th cent.</AUTHOR>
        <RESFSTMT><RESFSTMT creation of machine-readable version</RESFSTMT>
        <NAME>L'espion, Scrgo</NAME>
      </RESFSTMT>
      <TITLESTMT>
        <EXTENT>
          <SEG TYPE="designation">Text data</SEG>
          <SEG TYPE="size">(1 file : ca. 206 kilobytes)</SEG>
        </EXTENT>
        <PUBLICATIONSTMT><AUTHORITY ID="LISSE">Deposited by
          <NAME TYPE="PERSON">Laiguan, Sergio</NAME>
          <NAME TYPE="ORG">Inst. d'Etudes Mediecaires, Universite de Montreal</NAME>
          <ADDRESS><ADDRESS>Inst. d'Etudes Mediecaires, van</ADDRESS>
          <ADDRESSLINE>Universite de Montreal</ADDRESSLINE>
          <ADDRESSLINE>CP 6128 Succ A</ADDRESSLINE>
          <ADDRESSLINE>HEC 317 Montreal PQ</ADDRESSLINE>
          <ADDRESSLINE>Canada</ADDRESSLINE>
          <ADDRESS>
            <DATE><DATE>
          </ADDRESS>
        </PUBLICATIONSTMT>
        <DISTRIBUTOR>
          <NAME KEY="url" TYPE="organisation">Oxford Text Archive</NAME>
          <NAME TYPE="place">Oxford</NAME>
        </DISTRIBUTOR>
        <ADDRESS>
          <NAME KEY="name" TYPE="organisation">Oxford University Computing Services</NAME>
          <ADDRESSLINE>
          <ADDRESSLINE>13 Banbury Road</ADDRESSLINE>
          <ADDRESSLINE>Oxford</ADDRESSLINE>
        </ADDRESS>
      </TEIHEADER>

```

Martin Wynne, Opening the Archives

20

Conversion of OTA metadata (2)

```

<?xml version="1.0" encoding="UTF-8" ?>
<!DOCTYPE OLAC-Repository SYSTEM "olacrep.dtd">
<OLAC-Repository>
  <identity>
    <repositoryIdentifier>otarest</repositoryIdentifier>
    <repositoryName>Oxford Text Archive</repositoryName>
    <adminEmail>martin.wynne@ota.ahds.ac.uk</adminEmail>
    <collectionDescription>A test set of metadata records from the Oxford Text Archive
  </collectionDescription>
  </identity>
  <sets>
    <set>
      <setSpec></setSpec>
      <setName></setName>
    </set>
    <set>
      <setSpec></setSpec>
      <setName></setName>
    </set>
    <record>
      <header>
        <metadata>
          <dc:lang>"en">
            <contributor></contributor>
            <coverage>(1 file : ca. 97.6 kilobytes)</coverage>
            <creator></creator>
            <date></date>
            <description></description>
            <format>epo</format>
            <format.encoding></format.encoding>
            <format.markup></format.markup>
          </dc:lang>
        </metadata>
      </header>
    </record>
  </sets>
  </OLAC-Repository>

```

Martin Wynne, Opening the Archives

21

Validation

- XMetaL (or any XML parser) to check for well-formedness and conformance to the DTD (olacrep.dtd)
- <http://www.w3.org/2001/03/webdata/xsv?docAddr=http%3A%2F%2Fusers.ox.ac.uk%2F~martinw%2Folac%2Fota.xml&warnigs=on&keepGoing=on&style=offline> to check against XML Schema

Martin Wynne, Opening the Archives

22

Harvest

- Done manually by Steven Bird and Eva Banik from the XML file I supplied;
- In principle can be harvested from a perl script at the archive website which reads the metadata and looks for updates.

Martin Wynne, Opening the Archives

23

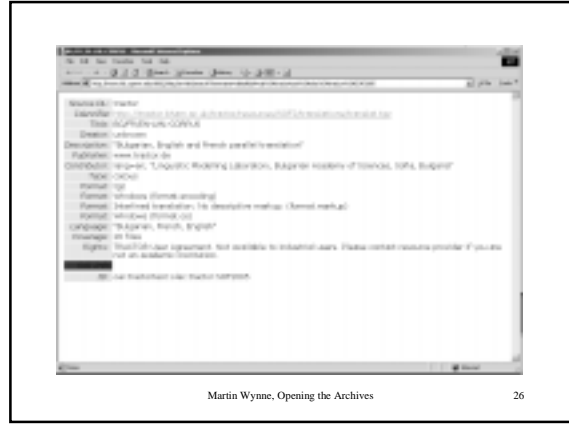


Martin Wynne, Opening the Archives

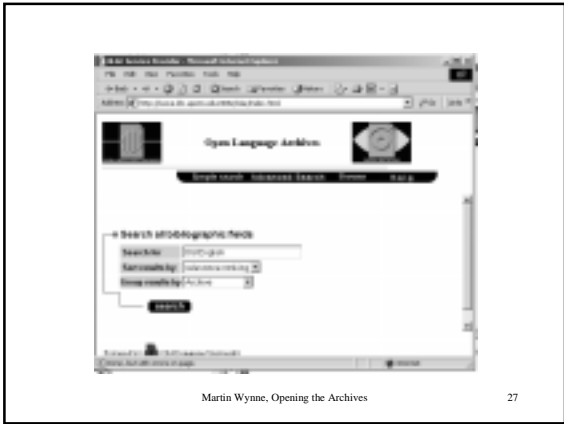
24



Martin Wynne, Opening the Archives



Martin Wynne, Opening the Archives



Martin Wynne, Opening the Archives



Martin Wynne, Opening the Archives



Martin Wynne, Opening the Archives



Results

- Metadata successfully harvested, and the virtual data provider accesses the records
- Recall imperfect
- Metadata imperfect
- Migration tools imperfect
- But the future is bright!

Acknowledgements

Thanks for their help to:

- Steven Bird (OLAC and LDC) plus Eva Banik and the ARC people (at Old Dominion University) as developers of the OLAC prototype
- TRACTOR people in Birmingham (esp. Everita Milconoka, Andrius Utka, Wolfgang Teubert)
- OTA people (John Leedham, Michael Popham and Alan Morrison)

Archives and Organisations

- <http://www.tractor.de>
- <http://ota.ahds.ac.uk>
- <http://www.open-archives.org>
- <http://www.language-archives.org>
- <http://www.mpi.nl/ISLE>
- <http://www.wave ldc.upenn.edu:8082/olac/>
- <http://www.linguistlist.org>

Documentation

- Warner, S, *Exposing and Harvesting Metadata Using the OAI Metadata Harvesting Protocol: A Tutorial*, <http://library.cern.ch/HEPLW/4/papers/3/>
- Simons & Bird (2001), *OLAC Metadata Set*, <http://www.language-archives.org/OLAC/olacms-20011022.html>
- Bird & Simons (2001), *The OLAC metadata set and controlled vocabularies. ACL Workshop on sharing tools and resources for research and education*, <http://arXiv.org/abs/cs/0105030>