# DSpace and Web Material: Inroads and Challenges

Leslie Myrick, NYU

DLF Spring Forum

April 15, 2005

# What I'll Be Covering

- NDIIPP "Web at Risk" Project
- Web Archive Data Object Modeling
- DSpace and HTML
  - Issues, Challenges, Prototypes
- Dspace, METS and Heritrix .arc format
- What's needed now ; desiderata for DSpace 2.0 +

# NDIIPP "Web at Risk" Partnership

- California Digital Library + UC partners
- University of North Texas
- New York University

- SDSC
- Stanford University
- Arizona State University
- Sun Microsystems

- Our LC AOTR is Martha Anderson

# Web Archiving 'Giants'

- Internet Archive Wayback Machine

- National Library of Australia PANDORA
- Royal Library of Sweden Kulturarw3

- LC MINERVA Project

NEW YORK UNIVERSITY

< > international internet preservation consortium - welcome - Microsoft Internet Explorer

File   Edit   View   Favorites   Tools   Help

Back   Search   Favorites   Media

Address  http://netpreserve.org/about/index.php   Go   Links »

# netpreserve.org
## international internet preservation consortium

contact
site search with google:
OK

**Welcome:**

about:
mission
members
working groups
contact
press releases

publications:
reports

software:
downloads

newsletter :

subscribe

## What's new?

**📅 July 20, 2004 / Publications / Report**
With the goal of constructing a live test bed for a Web crawling system, the Metrics and Testbed Working Group of the IIPC wrote a report which presents a taxonomy of challenges that a crawler may encounter on the Web at large when trying to copy content for Web archiving. more...

**📅 July 20, 2004 / Publications / Report**
The Metrics and Testbed Working Group of the IIPC conducted a survey which is an attempt to identify and classify many of the general conditions found on Web sites that influence the harvesting of content and the quality of an archival crawl. more...

**📅 May 5, 2004 / Press release**
In acknowledgement of the importance of international collaboration for preserving internet content for future generations, the International Internet Preservation Consortium was formed in 2003. more...

## About the consortium:

The national libraries of Australia, Canada, Denmark, Finland, France, Iceland, Italy, Norway, Sweden, The British Library (UK), The Library of Congress (USA) and the Internet Archive (USA) acknowledged the importance of international collaboration for preserving Internet content for future generations and therefore decided to form a consortium called the International Internet Preservation Consortium.

The goals of the consortium are:

- To enable the collection of a rich body of Internet content from around the world to be preserved in a way that it can be archived, secured and accessed over time.
- To foster the development and use of common tools, techniques and standards that enable the creation of

start    DSpace at New York ...    < > international inte...    Democratic Socialist ...    Jasc Paint Shop Pro    Microsoft PowerPoint ...    8:44 AM

# IIPC

- International Internet Preservation Consortium
  - National Library of Italy (Firenze)
  - Royal Library of Denmark
  - National Library of Finland
  - Internet Archive
  - Royal Library of Sweden
  - National Library of Iceland
  - Library and Archives / Canada
  - National Library of Norway
  - National Library of Australia
  - British Library
  - Library of Congress

# Local Expertise/Groundwork

- CDL
  - Web-based Government Information Project
  - California Recall Election Project
- UNT
  - CyberCemetery
- NYU
  - CRL Political Communications Web Archiving Project (Cornell, NYU, Stanford, UT Austin)
- LOC
  - MINERVA
  - IIPC Partnership

# Partnership's Four Paths

- Content Identification, Selection and Acquisition
  - Selection, Curatorial Issues
- Content Harvest and Analysis
  - Crawler
- Content Ingest, Retention and Transfer
  - CDL Digital Preservation Repository + DSpace
  - SRB, other grid technology, other means for synchronized replication
- Partnership Building
  - Technical and Human infrastructures

# What are we Capturing?

- WADO
  - Web Archive Digital Object
- Crawls, websites?
  - Many seed URLs recursively processed
    - archived websites, however they are defined
- Storage and Archive Format
  - Website mirrors
  - Flat hierarchy: entry page + all files
  - Gzipped archive files (.arc)

# Web Capture Tools

- Precrawl / Analysis Tools
  - Some of Linklint's functionality? And more
  - http://www.linklint.org/
    - Analysis of file types; "skipped" actions; missing pages
- Curatorial Interface
  - NLA's PANDAS as one model
  - Andy Boyko's PreCrawl tool for LC / IIPC
  - UIUC/OCLC Echo DEPository tools

# PANDAS Interface

# Crawler

- Crawler candidates
  - HTTrack
  - Heritrix

# HTTrack

- Developer: Xavier Roche & co.
- http://www.httrack.com/
- Written in C, open source
- Archive format: website mirror; .arc possible
- Core of NLA's PANDAS application
  - With Java UI modules

# Strengths and Weaknesses

- + Configurability
- + Small footprint
- + Incremental crawls possible
- - Scaling issues
- - No multi-machine crawling capability

# Heritrix

- Developers: Internet Archive et al
- http://crawler.archive.org/
- Working closely with IIPC
- Java, open-source ; Sourceforge
- Archive format: .arc.gz or mirror
- 1.0 release Aug 2004; up to 1.2.0 (Nov 2004); 1.4.0 (scheduled March 2005)

# Heritrix Architecture

- Scope – URIs to process
- Frontier -- controls processing
- Processor chains
  - Prefetch
  - Fetch
  - Extraction
  - Write
  - Postprocess

# Strengths and Weaknesses

- + International Standards
- + Programming base / experience
- + Configurability, extensibility
- - Too-sophisticated UI ; reports
- - Memory management
- - No multi-machine processing as yet
- - No incremental harvest capability as yet

# Heritrix .arc format (2.x)

- File header for .arc file itself

- Crawl metadata with arc and dc namespacing
  - Crawl host, operator, date

- Each harvested file and its metadata
  - [ DNS head]
  - HTTP response headers
  - The captured file

# The problem with .arcs

- Many seed URLs per crawl = many sites per .arc
- 100 MB limit (default) = many .arc files per crawl
- Management nightmare?
- Libarc tools promise to mitigate these I/O problems:
  - http://sourceforge.net/projects/libarc/

# The other problem with .arcs

- Gzipped format
  - each file and metadata bytestream is a separate member
- Solutions?
  - Interface with Heritrix API
    - can navigate a gzipped .arc:
      - ArcReader methods
      - Create a mini-Wayback machine
  - Use METS to manage transform object vs data object
- Question: How can METS and DSpace be tweaked to handle a zipped web archive?

# How does DSpace handle HTML?

# DSpace Data Model

# Objects in DSpace Data Model

- Item

  A technical report; a data set with accompanying description; a video recording of a lecture

- Bundle

  A group of HTML and image bitstreams making up an HTML document

- Bitstream

  A single HTML file; a single image file; a source code file

# DSpace and HTML: Webpage as Item

- Webpage as Item with its own handle
  - a bundle of files with a nominated primary bitstream (HTML)
- Integrated view of HTML page
  - together with all other files necessary for rendering (css, images, javascript).

# HTML Bundle before …

# ... HTML Bundle after

# Webpage-Level Item + and -

- + Search returns hits at page-level
- - Can't navigate natively within DSpace to any other webpage items from the same site that might be in the repository.
  – Could use METS DIP to navigate handles
- - Requires iterative loads of same files from other pages in same website (css, icons used across a site)

# Website as Item

- Website as Item
  - bundle of files with nominated primary bitstream (entry page's HTML).
- All other files (HTML and non-HTML) flattened out beneath this primary bitstream.

# Website-Level Item + and -

- + Files ingested and stored once
- + Can navigate hyperlink structure natively to DSpace
  - with some adjustments to archived HTML files.
- - Search returns hits at website level
- - HTML files' hyperlinks require some tweaking before loading

# DSpace Restrictions on HTML Objects

- No dynamic content (e.g. PHP, CGI)
- * All files must have unique names
- All hyperlinks must be relative and not refer to parents
  - myfile.html is okay
  - /myfile.html is not okay
  - Originally document-relative paths were anathema (../myfile.html)

# All Filenames Must be Unique

- Problematic in "real world" archiving
  - index.html at every directory level in site
- Use long filenames = path + filename

<snip>

paper/2003/feb-mar/SDFebMar03.rtf

paper/index.htm

statements/index.html

statements/nlc.html

students/index.html

students/war.html

index.html

</snip>

# Sidebar: Websites in DSpace:CWSpace

- William Reilly and Rob Wolfe (kudos and thanks)
- http://cwspace.mit.edu/
- Archiving MIT Open Courseware in DSpace
- IMS-CP packaging; conversion to METS AIP
- Impetus for changes to HTMLServlet and elsewhere; recent patch (11435750)
  - Larry Stone

# Recent Servlet Revisions

- Slashes in filenames inherently problematic
- Long filenames in contents file for ItemImporter can now be:
  - Handled by bitstreamServlet
  - Rendered by HTMLServlet as distinct items:

  \<snip\>

  paper/index.htm

  statements/index.html

  students/index.html

  index.html

  \</snip\>

# Adjustments to HTML files

- Document-relative paths within the files' hyperlinks now mandatory below root level (../)

- Necessary to navigate relative to what HTTPServlet perceives as "server root":
  - "server root" = website item handle
  - Concatenate long filename to it:
  - http://dspace.myu.edu/12345678/9/students/index.html

# DSpace and Heritrix .arc.gz: Problem and Solutions?

- Problem:  .arc.gz can be stored but not accessed as such; invokes gunzip, stuffit etc.

- Solutions?
  - Unzipped .arc could be manipulated locally by third party XSLT, e.g. METS Viewer
  - Interface with Heritrix API or with other tools that can handle .arc.gz format, e.g. ArcReader methods

- Caveat
  - Have to manage transform object over against data object

# DataObject vs transformObject

- XFDU dataObject vs transformObject
  - Transform nested in the dataObject

# XFDU DataObject and TransformObject

# Basic METS Website Data Model

<div> website
    <div> HTML page
        <fptr>
            <par>
                <area> bitstream (HTML)
                <area> bitstream (IMG)
            </par>
        </fptr>
        <div/> hyperlink on HTML page
    </div>
    …
</div>

# METS structMap for webpage

```
<METS:div DMDID="DM1" TYPE="web page" ID="page18"
    LABEL="www.apgawomen.org/index.html">
    <METS:fptr>
        <METS:par>
            <METS:area FILEID="FID18"/>                    [index.html ]
            <METS:area FILEID="FID1036"/>                  [notjust.swf]
            <METS:area FILEID="FID1043"/>                  [apgawnew.swf]
            <METS:area FILEID="FID1075"/>                  [enterarrow.gif]
        </METS:par>
    </METS:fptr>
    <METS:div TYPE="hyperlink" ID="LINK1" LABEL="home">
        <METS:fptr>
            <METS:area BEGIN="000" BETYPE="BYTE" END="111"
                FILEID="FID18"/>
        </METS:fptr>
    </METS:div>
```

# Mapping Hyperlink Structure

```
<METS:div DMDID="DM1" TYPE="web page" ID="page18" LABEL="www.apgawomen.org/index.html">
      <METS:fptr>
            <METS:par>
                <METS:area FILEID="FID18"/>                    [index.html ]
                <METS:area FILEID="FID1036"/>                  [notjust.swf]
                <METS:area FILEID="FID1043"/>                  [apgawnew.swf]
                <METS:area FILEID="FID1075"/>                  [enterarrow.gif]
            </METS:par>
      </METS:fptr>
       <METS:div TYPE="hyperlink" ID="LINK1" LABEL="home">
           <METS:fptr>
               <METS:area BEGIN="000" BETYPE="BYTE" END="111" FILEID="FID18"/>
           </METS:fptr>
         </METS:div>
</METS:div>

. . .

<METS:structLink>
   <METS:smLink from="LINK1" to="page1059" xlink:title="home"/>
   <METS:smLink from="LINK2" to="page113" xlink:title="officers"/>
   <METS:smLink from="LINK3" to="page102" xlink:title="calendar"/>
</METS:structLink>
```

# Possible adjustments to METS for .arc.gz

- Nested <file>s in <fileSec>
- <stream>s to handle metadata headers

- Possibly used along with either:
  - Nested transformObject for each file or
  - Dual structMaps
    - One for transform object
      - SM01: gzip object as the root div; .arc as child; all files as grandchildren;
    - One for website object unzipped
      - SM02: standard logical representation of files as they would exist outside of .arc.gz

# One possibility: Nested <file>s

```
<METS:fileSec>
   <METS:fileGrp>
      <METS:file ID="FID1" MIMETYPE="application/x-gzip" ADMID="ADM001">
         <METS:FLocat LOCTYPE="OTHER" OTHERLOCTYPE= "TRANSFORM-URI"
xlink:href="file:///usr/local/heritrix/jobs/Test07/arcs/IAH-20050203191213-00000-
euterpe.bobst.nyu.edu.arc.gz"></METS:FLocat>
               <METS:file ID="FID2" MIMETYPE="text/plain" ADMID="ADM01">
                  <METS:FLocat LOCTYPE="OTHER" OTHERTYPE="TRANSFORM-URI"
xlink:href="file://usr/local/heritrix/jobs/Test07/arcs/IAH-20050203191213-00000-
euterpe.bobst.nyu.edu.arc"></METS:FLocat>
               <METS:file ID="FID3" MIMETYPE=" text/html" ADMID="ADM1">
                  <METS:FLocat LOCTYPE="URL" xlink:href="www.apgawomen.org/"></METS:FLocat>
               </METS:file>
               [ other website members / files here ]
            </METS:file>
      </METS:file>
   </METS:fileGrp>
</METS:fileSec>
```

# METS StructMap01: Transform Object

```
<METS:structMap ID="SM01" TYPE="logical-transformation">
    <METS:div DMDID="DM01" TYPE="web site" ID="page1" LABEL="www.apgawomen.org/">
        <METS:fptr>
                <METS:par>
                    <METS:area BEGIN="1725" BETYPE="BYTE" FILEID="FID1"></METS:area>
                    <METS:area BEGIN="6571" BETYPE="BYTE" FILEID="FID1"></METS:area>
                    <METS:area BEGIN="2670" BETYPE="BYTE" FILEID="FID1"></METS:area>
                    <METS:area BEGIN="17561" BETYPE="BYTE" FILEID="FID1"></METS:area>
                </METS:par>
        </METS:fptr>
        <METS:div TYPE="hyperlink" ID="LINK1" LABEL="home">
            <METS:fptr>
                    <METS:area BEGIN="000" BETYPE="BYTE" FILEID="FID3"></METS:area>
                </METS:fptr>
        </METS:div>
            . . .
    </METS:div>
</METS:structMap>
```

# METS StructMap02
# Unzipped Website Object

```
<METS:structMap ID="SM02" TYPE="logical">
    <METS:div DMDID="DM01" TYPE="web site" ID="page18" LABEL="www.apgawomen.org/">
        <METS:fptr>
            <METS:par>
                <METS:area FILEID="FID3"></METS:area>
                    <METS:area FILEID="FID1036"></METS:area>
                    <METS:area FILEID="FID1043"></METS:area>
                    <METS:area FILEID="FID1075"></METS:area>
            </METS:par>
        </METS:fptr>
        <METS:div TYPE="hyperlink" ID="LINK1" LABEL="home">
            <METS:fptr>
    <METS:area BEGIN="000" BETYPE="BYTE" END="111" FILEID="FID3"></METS:area>
            </METS:fptr>
        </METS:div>
        . . .
    </METS:div>
</METS:structMap>
```

# DSpace Desiderata: Ingest

- Load Scripts (Website SIP Client):
  - Script to visit files in load directory and write contents page for ItemImport
  - Script to automate dublin_core.xml
  - Script to "correct" hyperlinks in HTML
    - All pages siblings below root with ../ path
    - Disable mailtos, cgis, etc.
    - Archive or correct links to external pages
- Functionality to automate nomination of primary bitstream
- METS Import
  - Finer-grained metadata possible

# DSpace Desiderata: Archival Storage

- Functionality to manage version control issues
  - Successive snapshots of same site
  - Migration of bitstreams
- Facilitate cross-collection item linking

# DSpace Desiderata: Metadata

- Ability to apply (and discover) more metadata at bitstream level
  - Controlled descriptive information for HTML, PDF, etc
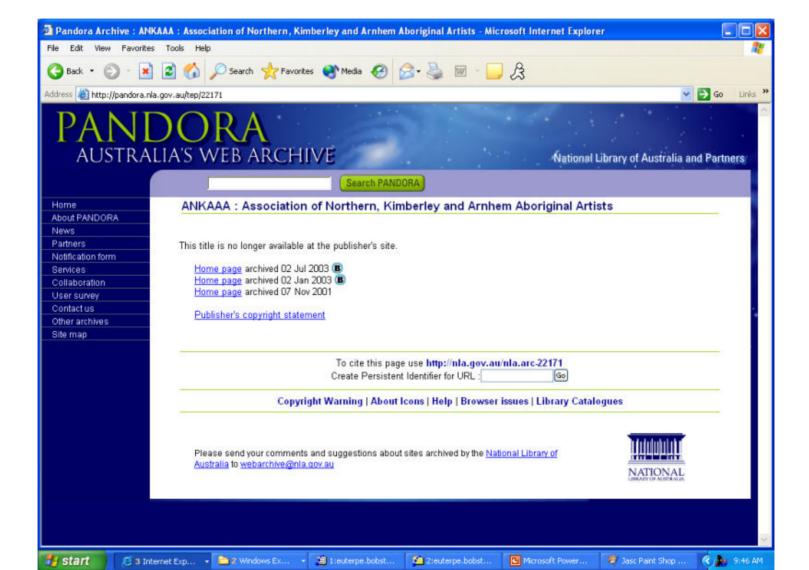  - Extend dublin_core.xml SIP to bitstream level?
  - Develop METS SIP?
- METS AIP
  - First of all, a structMap (or two)
  - Manage transform objects' (.arc, .gz) relationship to data objects (zipped files)
  - Reflect complexity of hyperlinked object
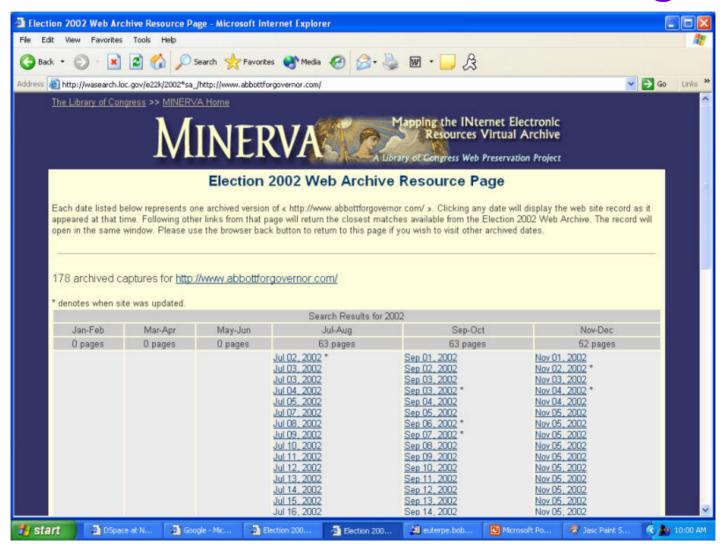    - structLink cross-reference to structMap

# DSpace Desiderata: Access

- Display search results at bitstream level
  - for HTML, PDF, MSWord, etc
- Manage navigation of different versions of same site
  - Title Entry Page (NLA PANDAS)
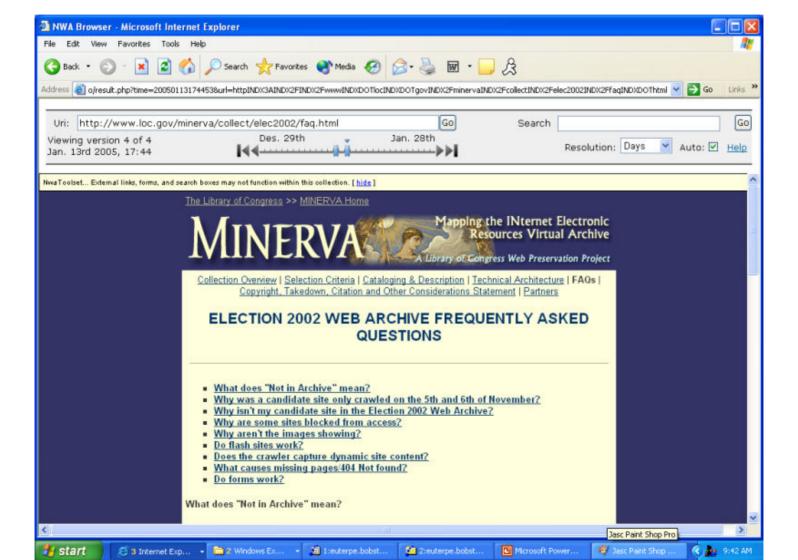  - IA Resource Page (Wayback, MINERVA)
  - Timeline (NWA Toolset)

# PANDAS TEP Interface

# IA /MINERVA Resource Page

# NWA Toolset Timeline Access

# DSpace < > Heritrix

- Use METS for supplemental management of transform vs unzipped object
- Interface with Heritrix API
  - To write SIPs from the .arc
  - For possible creation of mirrors
    - MirrorWriterProcessor
  - For Wayback functionality
    - ArcReader

# For More Information

- <u>leslie@nyu.edu</u>
- <u>http://dlibdev.nyu.edu/demo/DLF2005.html</u>
- DPPWAR/ NYU DSpace Testbed
  - <u>http://dlibdev.nyu.edu/dspace/handle/123456789/1</u>

Go to Demo