# Using Schematron for Analyzing Conformance to Best Practices for EAD, TEI, and MODS

## (and some other thoughts on workflow tools)

Jenn Riley

Metadata Librarian

Indiana University Digital Library Program

---

# Consistency is a challenge

- Document-centric XML (TEI, EAD) is very difficult to create consistently
- Some common tools to help:
  - Schema/DTD validation
  - Tag libraries
  - XML templates
  - Example documents
  - Keyboard macros
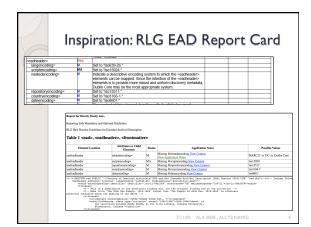  - Detailed encoding guidelines
- These are not enough!

---

# Another possible tool layer

- Machine validation of a file against local encoding guidelines
- Can only go so far, but that far is extremely helpful
- Indiana University implemented using:
  - Schematron (http://www.schematron.com)
  - <oXygen /> plugin architecture

---

# Inspiration: RLG EAD Report Card

---

---

## More info on Schematron

- ISO/IEC 19757 - Document Schema Definition Languages (DSDL) - Part 3: Rule-based validation – Schematron.
- Be careful! http://www.schematron.com has the primary specs; http://schematron.com is for a particular company's tool using them. (Weird.)
- This is the page you want:

## Using a Schematron file

- Schematron home page provides two distributions:
  - One for XSLT 1.0 processors and one for 2.0 processors
  - Each includes a set of three stylesheets to be used in turn on the Schematron file
  - Result of this processing is a stylesheet to be run on your XML instance document
- IU implementation wraps this all up into an <oXygen /> plugin written in Java
- You could also pipe them together with a shell script, a Windows .bat file, etc……

---

### Levels of Adoption

The Digital Library Federation / Aquifer Implementation Guidelines for Shareable MODS Records provide guidance on the use of M Aquifer initiative. The Guidelines, however, are a record-centric view of Aquifer's goals, whereas it is often helpful to set priorities for document, the Digital Library Federation MODS Guidelines Levels of Adoption, describes five general categories of user fu specific recommendations from the Guidelines. It attempts to provide additional guidance to MODS implementers in the planning p possible when certain elements of the Guidelines are followed.

The Working Group welcomes comments on this documentation. Comments and questions can be sent to any Working Group mer

The Levels of Adoption are listed here from loosest to strictest.

**1) Minimum for participation: Allows users to cite the resource**

The minimum for participation level defines the information necessary for the most basic indexing of records. It represents a slightly decision reflects the view that Aquifer participants are leaders in the digital library field and therefore will likely be able to commit to
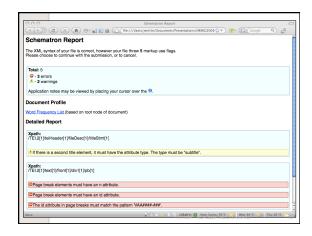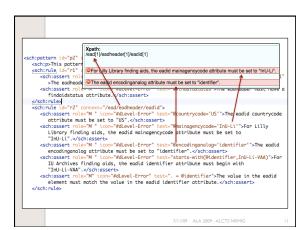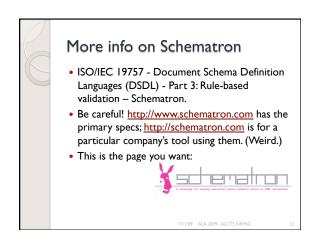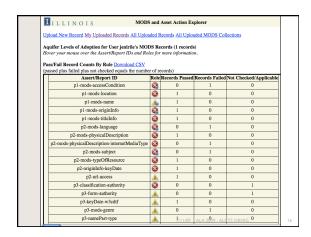
Aquifer will not harvest records that do not contain this information. This determination will be made at the set level.

- Syntactically valid MODS
- <name> if applicable (and known)
- at least one <titleInfo> element with one <title> subelement
- if known, at least one <originInfo> element with at least one <date> subelement from the recommended list
- one or more <location> elements with a <url> subelement, one of which should point to the resource in context
- <accessCondition> if anything other than completely unrestricted

**2) Minimum for doing anything useful: Allows users to perform basic searches and filtering**

This second level of adoption represents a minimum that will allow an institution's resources to be incorporated *meaningfully* into an minimum for participation level in order for the Aquifer project to obtain as much metadata as possible for the purposes of testing v believes that records that do not meet at least this second level of adoption are likely to be discoverable by end-users or usable in t time all Aquifer participants will meet at least this second level of adoption.

- <subject>
- at least one <language> element for all resources in which language is primary to understanding the resource
- One and only one date element must be marked as a key date
- At least one <typeOfResource> using values from the approved list
- Use of access attribute on location/url, or include only one URL that points to a view of the resource in context
- At least one <physicalDescription> element with <internetMediaType> subelement; value to be taken from appropriate list

---

### MODS and Asset Action Explorer

**I** ILLINOIS

Upload New Record  My Uploaded Records  All Uploaded Records

**MODS Collections Uploaded By Everyone**

[danc] Check Aquifer Levels of Adoption

(1 records) [3/21/2009 11:56:18 PM] View Check Aquifer Levels of Adoption
(0 records) [3/21/2009 9:31:23 PM] View Check Aquifer Levels of Adoption
(0 records) [3/21/2009 9:30:42 PM] View Check Aquifer Levels of Adoption
(0 records) [3/21/2009 9:12:16 PM] View Check Aquifer Levels of Adoption

**17964** [khage] Check Aquifer Levels of Adoption

(45 records) [1/16/2009 1:06:30 PM] View Check Aquifer Levels of Adoption

**acsc** [khage] Check Aquifer Levels of Adoption

(59 records) [1/15/2009 2:40:25 PM] View Check Aquifer Levels of Adoption

**ACSC** [thabing] Check Aquifer Levels of Adoption

(59 records) [1/12/2009 10:25:47 AM] View Check Aquifer Levels of Adoption

**AH** [khage] Check Aquifer Levels of Adoption

(34 records) [1/15/2009 2:38:51 PM] View Check Aquifer Levels of Adoption

---

### MODS and Asset Action Explorer

**I** ILLINOIS

Upload New Record  My Uploaded Records  All Uploaded Records  All Uploaded MODS Collections

**Aquifer Levels of Adoption for User jenlrile's MODS Records (1 records)**
*Hover your mouse over the Assert/Report IDs and Roles for more information.*

**Pass/Fail Record Counts By Rule** Download CSV
(passed plus failed plus not checked equals the number of records)

| Assert/Report ID | Role | Records Passed | Records Failed | Not Checked/Applicable |
|---|---|---|---|---|
| p1-mods-accessCondition | | 0 | 1 | 0 |
| p1-mods-location | | 1 | 0 | 0 |
| p1-mods-name | | 1 | 0 | 0 |
| p1-mods-originInfo | | 1 | 0 | 0 |
| p1-mods-titleInfo | | 1 | 0 | 0 |
| p2-mods-language | | 0 | 1 | 0 |
| p2-mods-physicalDescription | | 1 | 0 | 0 |
| p2-mods-physicalDescription-internetMediaType | | 0 | 1 | 0 |
| p2-mods-subject | | 0 | 1 | 0 |
| p2-mods-typeOfResource | | 1 | 0 | 0 |
| p2-originInfo-keyDate | | 1 | 0 | 0 |
| p2-url-access | | 1 | 0 | 0 |
| p3-classification-authority | | 0 | 0 | 1 |
| p3-form-authority | | 0 | 0 | 1 |
| p3-keyDate-w3cdtf | | 1 | 0 | 0 |
| p3-mods-genre | | 0 | 1 | 0 |
| p3-namePart-type | | 1 | 0 | 0 |

---

## Let's step back

- How can better tools revolutionize metadata creation workflows?
  - Promoting consistency
    - This is hard and not something that humans are generally good at
  - True interoperability between systems
    - *Without futzing!*
- We spend too much valuable human time doing repetitive and low-value tasks as part of descriptive workflows

---

## Were do we go from here?

- Make better use of available technologies
  - Automating
  - Streamlining
  - Validating
- We can and *must* do our jobs better and more efficiently, with the help of better tools
  - Providing comparable services with less
  - Creating a convincing argument for more?

## There is no excuse for not having usable metadata creation tools.

- Smart systems are possible and necessary
  - Configurable
  - Modular
  - Connected
- Make it easy to do it well
  - Consistent
  - Complete
  - Efficient
- Make it hard to do it poorly
- *We must pay attention to user interface design for cataloging tools*

## OK, rant over. Thank you!

- jenlrile@indiana.edu
  - (watch out for the invisible "l" in the middle)
- Slides and handout:
  - On ALA presentations Wiki <http://presentations.ala.org>
  - On my home page <http://www.dlib.indiana.edu/~jenlrile/presentations/nrmig2009/>