

# Practical Approaches for Migrating Metadata to Hyrax

Samvera Virtual Connect 2017  
July 18th, 11:55-12:20 EDT

The background of the slide is a photograph of a forest. The scene is filled with tall, dark evergreen trees. Sunlight filters through the canopy from above, creating bright highlights on the trunks and some low-hanging branches. The ground is covered with a mix of green grass and fallen brown leaves.

**Link to the slides, notes, links:  
<http://bit.ly/HyraxDataMigration>**

# **Today's Facilitators**

**Christina**

Data Ops Repository Specialist, Stanford Univ. Libraries,  
@cm\_harlow

**Tom**

Ph.D. Student at University of Wash. iSchool; DCE, @no\_reply

# This Session's Goals: Help You...

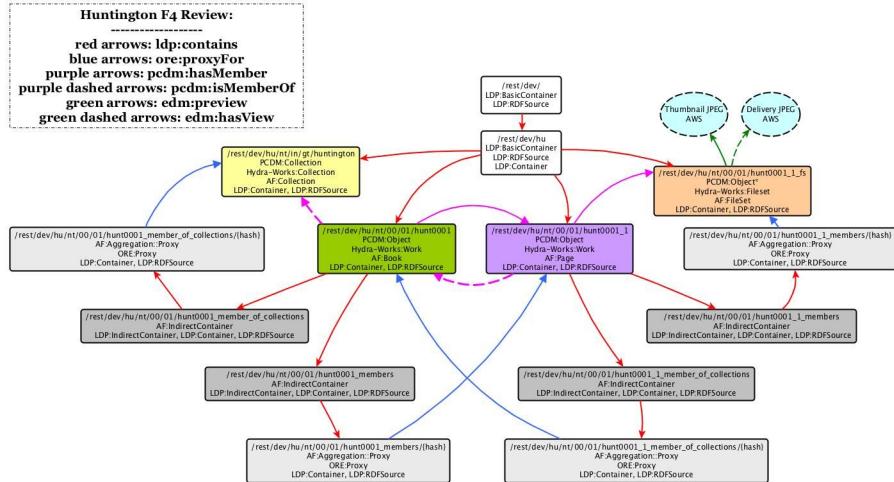
- Understand the Hyrax Metadata “Status Quo”
- Know Where to Start with Metadata Migration
- Point Out Questions, Priorities, & “Gotchas”
- Point You to Helpful Resources & Communities

The background of the slide features a photograph of a dense forest. The foreground is filled with the dark trunks and branches of tall evergreen trees. In the distance, more trees are visible through a layer of mist or fog, creating a sense of depth. The overall atmosphere is mysterious and natural.

# Status Quo of Hyrax Metadata

# Object Models & Application Profiles

## Object Model



## Application Profile

### PCDM:Collection > HydraWorks:Collection : Digital Collection

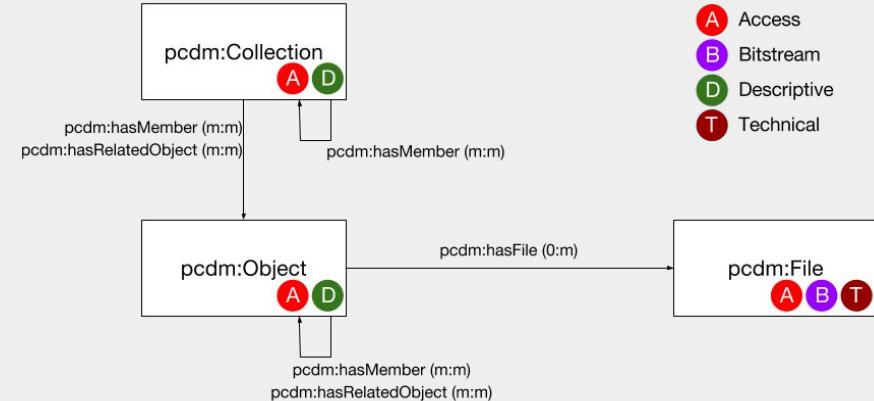
This is the collection resource representing Huntington Digital Collection.

#### Descriptive Profile

predicate	value	notes
dcterms:title	"Huntington Free Library Native American Collection"^^xsd:string	Need language typing
dcterms:abstract	"One of the largest collections of books and manuscripts of its kind, the Huntington collection contains extensive materials documenting the history, culture, languages, and arts of the native tribes of both North and South America. Contemporary politics and human rights issues are also important components of the collection. Full text of a selection of 91 books from the Huntington Free Library Native American Collection representing the various genres in the collection."^^xs:string	Need language typing
dcterms:date	"2010"^^ <a href="http://id.loc.gov/datatypes/edtf">http://id.loc.gov/datatypes/edtf</a>	Need date (data type) typing. Make sure is not dcterms:created
dcterms:identifier	"6790930"	identifier typing to be added in phase 2 of migration.
		To be removed

# PCDM & Hydra Works

Default Usage in Hyrax



- PCDM is a low-level object model for repository data.
- Hydra Works extends PCDM by breaking Objects into Works and FileSets.

*Why shared models?*

# PCDM:Collections MAP

Default Usage in Hyrax

Collection's members are Objects and/or other Collections.

## Hyrax::CoreMetadata

- Title (dct:title)
- Depositor (marcrel:dpt)
- Date Uploaded (dct:dateSubmitted)
- Date Modified (dct:modified)

## Hyrax::BasicMetadata

Creator, Contributor, Rights, Description, &c.

See:

[http://samvera.github.io/customize-metadata-mode\\_1.html#basic-metadata](http://samvera.github.io/customize-metadata-mode_1.html#basic-metadata)

# HW:Work MAP

Default Usage in Hyrax

‘A work or intellectual entity, such as a book, film, dissertation, etc.’

*A hw:Work is a pcdm:Object.*

Use Works for repository objects. I.e. *books*, but also *pages* (if you handle pages).

**Core & Basic Metadata**

# HW:FileSet MAP

Default Usage in Hyrax

“A group of related Files”

*A hw:FileSet is a pcdm:Object.*

Use FileSets for groupings of files that aren’t repository objects; e.g. a scanned text and related OCR output.

**Core & Basic Metadata**

Use FileSet metadata for information about the group.

FileSets may be members of more than one Object.

# PCDM:File MAP

Default Usage in Hyrax

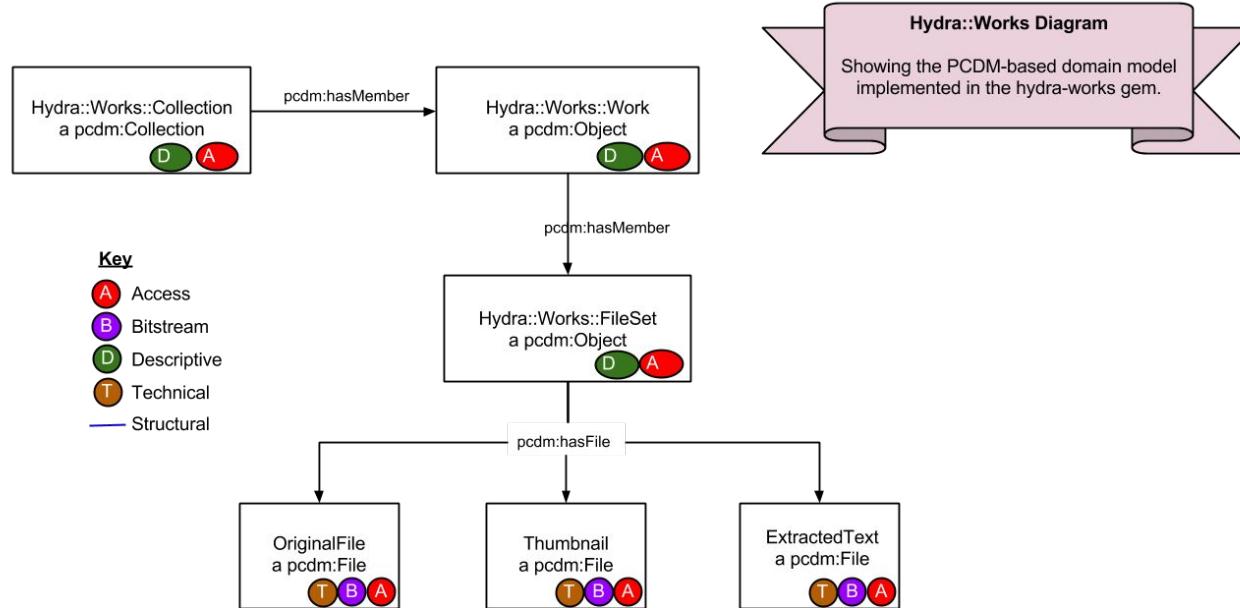
“A File is a sequence of binary data”

Files are members of **one** Object.

File Metadata:

- Label (rdfs:label)
- File Name (ebu:filename)
- File Size (ebu:filesize)
- Date Created (ebu:dateCreated)
- Date Modified (ebu:dateModified)
- Byte Order (sweetjpl:byteOrder)
- File Hash (premis:hasMessageDigest)

# Works Model



# Hyrax Basic Metadata

- Basic Metadata is intended to capture general use cases.
- Can be relatively easily changed (or not used at all) for Works.
- See:  
[http://samvera.github.io/  
customize-metadata-model.  
html#basic-metadata](http://samvera.github.io/customize-metadata-model.html#basic-metadata)

Property	Predicate	Multiple
label	ActiveFedora::RDF::Fcrepo::Model.downloadFilename	FALSE
relative_path	::RDF::URI.new('http://scholarsphere.psu.edu/ns#relativePath')	FALSE
import_url	::RDF::URI.new('http://scholarsphere.psu.edu/ns#importUrl')	FALSE
resource_type	::RDF::Vocab::DC.type	TRUE
creator	::RDF::Vocab::DC11.creator	TRUE
contributor	::RDF::Vocab::DC11.contributor	TRUE
description	::RDF::Vocab::DC11.description	TRUE
keyword	::RDF::Vocab::DC11.relation	TRUE
license	::RDF::Vocab::DC.rights	TRUE
rights_statement	::RDF::Vocab::EDM.rights	TRUE
publisher	::RDF::Vocab::DC11.publisher	TRUE
date_created	::RDF::Vocab::DC.created	TRUE
subject	::RDF::Vocab::DC11.subject	TRUE
language	::RDF::Vocab::DC11.language	TRUE
identifier	::RDF::Vocab::DC.identifier	TRUE
based_near	::RDF::Vocab::FOAF.based_near	TRUE
related_url	::RDF::RDFS.seeAlso	TRUE
bibliographic_citation	::RDF::Vocab::DC.bibliographicCitation	TRUE
source	::RDF::Vocab::DC.source	TRUE

# Adding a Field

- Select a predicate
- Select an accessor name
- Single- or Multi-value?
- Controlled Vocabulary?
- Required? Other validations?
- Ordered? (are you really modeling order, or something else?)
- Data types (e.g. dates, integers) and language tags are handled natively.

## Documentation -

- Defining Metadata  
<http://samvera.github.io/customize-metadata-model.html>
- Controlled Vocabularies  
<http://samvera.github.io/customize-metadata-controlled-vocabulary.html>

## Underlying Model -

- Resource Description Format (RDF)  
<https://www.w3.org/TR/rdf11-concepts/>

```
# Generated via
#   `rails generate hyrax:work GenericWork`
class GenericWork < ActiveFedora::Base
  include ::Hyrax::WorkBehavior
  include ::Hyrax::BasicMetadata
  # Change this to restrict which works can be added as a child.
  # self.valid_child_concerns = []
  validates :title, presence: { message: 'Your work must have a title.' }

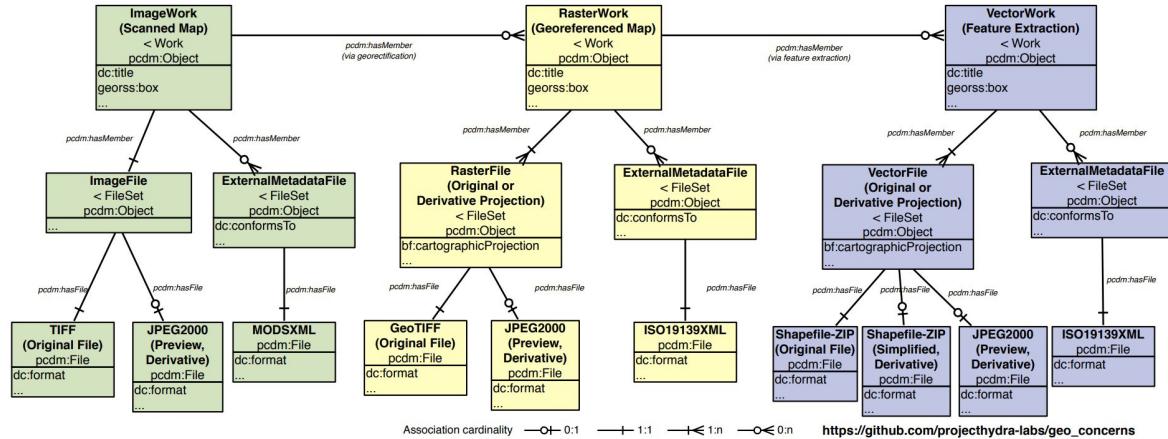
  property :contact_email, predicate: ::RDF::Vocab::VCARD.hasEmail, multiple: false do |index|
    index.as :stored_searchable
  end

  property :contact_phone, predicate: ::RDF::Vocab::VCARD.hasTelephone do |index|
    index.as :stored_searchable
  end

  property :department, predicate: ::RDF::URI.new("http://lib.my.edu/departments"), multiple: false do |index|
    index.as :stored_searchable, :facetable
  end
end
```

# Build, Reuse & Share MAPs

- Share both structure and application profiles.
- Spread development and maintenance labor around.
- Reusing saves work!

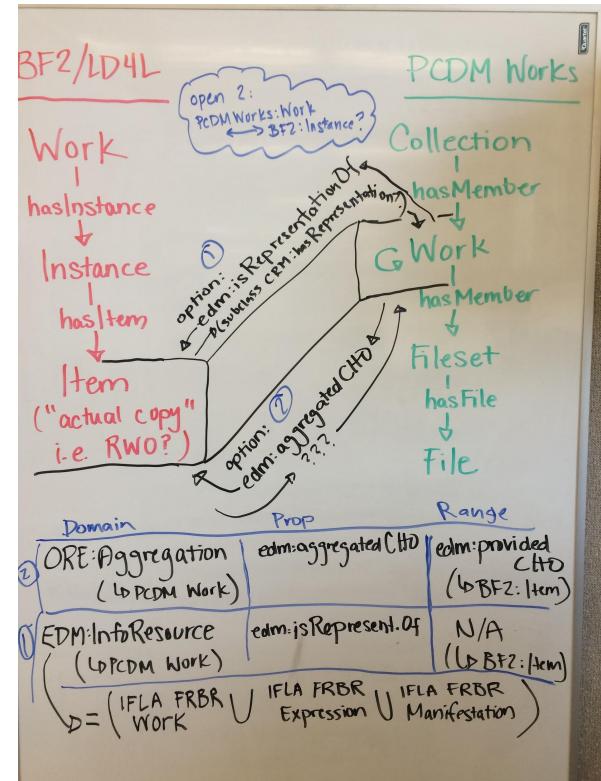


A photograph of a forest fire. In the foreground, there are fallen logs and green moss-covered rocks. The background shows tall, thin trees with white bark and orange flames. A large amount of white smoke is rising from the ground and between the trees.

# Where to Start in Migration Planning

# Questions to Work Out

- Pre-Migration  
    ● Classes of Objects Described? Types of Objects?
- Relationships between Object Classes? Nested Objects? “Complex” Objects?
- 1st Class Resources? 2nd Class? Context Classes?
- External Data? Entity Resolution? Look-Ups? New URIs? Identifiers?
- Discovery Access Points & Facets?
- Who Adds What Metadata?
- Metadata Current State? Fix What Errors When?
- Separate out Migration & Enhancements And...
- **DO NOT AIM FOR PERFECT**



# Performing Assessment

Do Your Metadata Homework,  
Using Whatever Tools Work For  
You.

```
400 b'Fact sheet'
  1 b'Fact sheet and poster'
  1 b'Fact sheet and posters'
3875 b'Image'
  1 b'Instruction Manual'
  45 b'Interview'
125 b'Learning Object'
  2 b'Literature Review'
  3 b'Musical Score'
275 b'Newsletter'
1270 b'Other'
  243 b'Pamphlets'
5224 b'Paper or Project'
  20 b'Paper or project'
6990 b'Periodical'
  3 b'Poster'
  37 b'Preprint'
759 b'Presentation'
  1 b'Proceedings'
  59 b'Project Plan'
338 b'Project Report'
119 b'Report'
  31 b'Software'
  426 b'Sound'
229 b'Syllabus'
3330 b'Technical Report'
  19 b'Technical report'
  1 b'Video/Moving Image and Poster'
721 b'Video/Moving Image'
  9 b'Video/Moving Object'
33 b'Website'
  6 b'Working Paper'
535 b'article'
600 b'book chapter'
  39 b'book'
  1 b'case study'
  1 b'conference proceedings'
  3 b'dataset'
6158 b'dissertation or thesis'
  4 b'fact sheet'
  1 b'final report'
13 b'learning object'
  1 b'literature review'
3873 b'newsletter'
  21 b'other'
  4 b'paper or project'
  4 b'periodical'
  1 b'preprint'
73 b'presentation'
  1 b'procedure manual'
  1 b'project research'
932 b'report'
  53 b'syllabus'
  9 b'technical report'
11 b'video/moving image'
  1 b'video/moving image'
  ↵<sul.cmlarlow @ li-dl-2628-db6b in ~/C/s/g/c/metadataQA>
-><(master)* >-> exit
-><sul.cmlarlow @ li-dl-2628-db6b in ~/C/s/g/c/metadataQA>
-><(master)* >-> ttypif myrecording
```

# Clarify & Iterate on Expectations

Assess, Prioritize  
Functionalities, Normalize,  
Migrate, Normalize, Enhance

The RDF relationships expected for this class, symmetry of those properties not assumed (hence A -> B and B -> A are both given)

Structure of Digital Objects in our Fedora 4/Hydra stack using PCDM:

**PCDM:Collection > HydraWorks:Collection : Digital Collection**

This is the digital collection that current maps to the dlxs identifier sets (i.e. 'hunt', 'bol', etc.). This is required. Can be repeated/nested as needed.

Descriptive Metadata Profile:

- dcterms:abstract = string
- dcterms:alternative = string
- ~~Irdau:P60376 [curator]~~ = string (phase 2 removal)
- ~~Irdau:P60376 [curator]~~ = Agent URI (phase 2 addition)
- dcterms:date = EDTF literal
- dcterms:identifier = string
- ~~dc:language~~ = string (phase 2 removal)
- ~~dc:publisher~~ = string (phase 2 removal)
- ~~dcterms:publisher~~ = Agent URI (phase 2 addition)
- dcterms:relation = URL
- ~~dc:subject~~ = string (phase 2 removal)
- ~~dcterms:subject~~ = Concept URI (phase 2 addition)
- dcterms:title = literal

Structural Profile:

If no secondary PCDM:Collection for Set:

- *Digital Collection -PCDM:hasMember-> Intellectual Work*
- *Digital Collection <-PCDM:isMemberOf- Intellectual Work*

If there is a secondary PCDM:Collection for Set:

- *Digital Collection -PCDM:hasMember-> Set*
- *Digital Collection <-PCDM:isMemberOf- Set*

# “Graceful Degradation”

What is the shared,  
underlying logic that your  
metadata can revert to?

Full Representation

CIDOC-CRM    FRBR-OO    BIBFRAME

EDM    VRACore-RDF    MODSRDF

DPLA              Schema

DC Profiles

Common Abstractions

A photograph of a forest scene. The foreground is filled with tall, thin trees, likely pines or similar conifers, their trunks dark and straight. Sunlight filters down from the canopy above, creating bright, glowing vertical streaks on the tree trunks and illuminating patches of green grass and fallen leaves on the forest floor. The background is a dense wall of trees, their branches reaching out towards the center.

# Some Rough / Common Groups of Recommendations

# DSpace

- OAI-PMH DIMS Feed is Most Complete
- Beware Namespaces
- Complex Fields via Elements & Qualifiers
- Break Out Object Versions, File Info, Identifiers for Parts / Filesets / Objects
- Beware Language Tags

```
400 b'Fact sheet'
  1 b'Fact sheet and poster'
  1 b'Fact sheet and posters'
3875 b'Image'
  1 b'Instruction Manual'
  45 b'Interview'
125 b'Learning Object'
  2 b'Literature Review'
  3 b'Musical Score'
275 b'Newsletter'
1270 b'Other'
  243 b'Pamphlets'
5224 b'Paper or Project'
  20 b'Paper or project'
6990 b'Periodical'
  3 b'Poster'
  37 b'Preprint'
759 b'Presentation'
  1 b'Proceedings'
  59 b'Project Plan'
338 b'Project Report'
119 b'Report'
  31 b'Software'
426 b'Sound'
229 b'Syllabus'
3330 b'Technical Report'
  19 b'Technical report'
  1 b'Video/Moving Image and Poster'
721 b'Video/Moving Image'
  9 b'Video/Moving Object'
33 b'Website'
  6 b'Working Paper'
535 b'article'
600 b'book chapter'
  39 b'book'
  1 b'case study'
  1 b'conference proceedings'
  3 b'dataset'
6158 b'dissertation or thesis'
  4 b'fact sheet'
  1 b'final report'
  13 b'learning object'
  1 b'literature review'
3873 b'newsletter'
  21 b'other'
  4 b'paper or project'
  4 b'periodical'
  1 b'preprint'
  73 b'presentation'
  1 b'procedure manual'
  1 b'project research'
932 b'report'
  53 b'syllabus'
  9 b'technical report'
  11 b'video/moving image'
  1 b'video/moving image'
  ↵sul.cmlharlow@li-dl-2628-db6b in ~/C/s/g/c/metadataQA>
↳ (master)* ➤➤➤ exit
↳ sul.cmlharlow@li-dl-2628-db6b in ~/C/s/g/c/metadataQA>
↳ (master)* ➤➤➤ ttyp1f myrecording
```

# Bepress (Digital Commons)

- OAI-PMH Metadata-Export Is Most Complete
- Check Completeness of Data Dumps
- Some Regular Fields
- Authors Broken Into Multi-Parts
- Local Fields Beware
- Supplemental Files => Filesets, Related Objects, Parts

	abstract:  =====	8318/9895   84%
Agents	articleid:  =====	9895/9895   100%
Agents	authors/author/email:  =====	2270/9895   22%
Agents	authors/author/fname:  =====	5315/9895   53%
Agents	authors/author/institution:  =====	3874/9895   39%
Agents	authors/author/lname:  =====	5335/9895   53%
Agents	authors/author/mname:  =====	2131/9895   21%
Agents	authors/author/organization:  =====	3243/9895   32%
Agents	context-key:  =====	9895/9895   100%
Agents	coverpage-url:  =====	9895/9895   100%
Agents	disciplines/discipline:  =====	5043/9895   50%
Agents	document-type:  =====	9895/9895   100%
Agents	fields/field/value:  =====	9089/9895   91%
Agents	fields/field[@name=academic_advisor][@type:string]:  =====	30/9895   0%
Agents	fields/field[@name=academic_email][@type:string]:  =====	30/9895   0%
Agents	fields/field[@name=acknowledgements][@type:string]:  =====	1/9895   0%
Agents	fields/field[@name=advisor1][@type:string]:  =====	1576/9895   15%
Agents	fields/field[@name=advisor1_email][@type:string]:  =====	32/9895   0%
Agents	fields/field[@name=affiliation][@type:string]:  =====	1915/9895   19%
Agents	fields/field[@name=custom_citation][@type:string]:  =====	1123/9895   11%
Agents	fields/field[@name=degree_name][@type:string]:  =====	2047/9895   20%
Agents	fields/field[@name=department][@type:string]:  =====	1975/9895   19%
Agents	fields/field[@name=distribution_license][@type:string]:  =====	20/9895   0%
Agents	fields/field[@name=embargo_date][@type:date]:  =====	5483/9895   55%
Agents	fields/field[@name=geolocate][@type=special]:  =====	50/9895   0%
Agents	fields/field[@name=identifier][@type:string]:  =====	31/9895   0%
Agents	fields/field[@name=issnum][@type:string]:  =====	19/9895   0%
Agents	fields/field[@name=language][@type:string]:  =====	7764/9895   78%
Agents	fields/field[@name=latitude][@type:string]:  =====	48/9895   0%
Agents	fields/field[@name=link_institution][@type:string]:  =====	21/9895   0%
Agents	fields/field[@name=longitude][@type:string]:  =====	48/9895   0%
Agents	fields/field[@name=publication_date][@type=date]:  =====	9028/9895   91%
Agents	fields/field[@name=publication_date_date_format][@type:string]:  =====	226/9895   2%
Agents	fulltext-url:  =====	8006/9895   80%
Agents	keywords/keyword:  =====	6609/9895   66%
Agents	label:  =====	9895/9895   100%
Agents	native-url:  =====	2196/9895   22%
Agents	publication-date:  =====	9895/9895   100%
Agents	publication-title:  =====	9895/9895   100%
Agents	submission-date:  =====	9895/9895   100%
Agents	submission-path:  =====	9895/9895   100%
Migration Aid	supplemental-files/file/archive-name:  =====	98/9895   0%
Files	supplemental-files/file/description:  =====	30/9895   0%
Parts	supplemental-files/file/mime-type:  =====	98/9895   0%
Parts	supplemental-files/file/upload-name:  =====	98/9895   0%
Parts	supplemental-files/file/url:  =====	98/9895   0%
	title:  =====	9895/9895   100%
	type:  =====	9895/9895   100%

# ContentDM

- Check TSV Export First for Most Complete
- Watch Out For Page-Level Objects Mapped to PCDM Objects / “Parts”
- Split Concatenated Fields
- Beware Local Fields

Refine OPEN oai php Permalink

Facet / Filter Undo / Redo 1

Using facets and filters

Use facets and filters to select subsets of your data to act on. Choose facet and filter methods from the menus at the top of each data column.

Not sure how to get started?  
[Watch these screencasts](#)

200 records

Show as: rows records Show: 5 10 25 50 records

record - metadata	record - metadata - oai_qdc:qualifieddc - dc:identifier	record - metadata	record - metadata	record - metadata	record - metadata
Architecture; Architects--Tennessee--Nashville; Bridges--Tennessee--Goodlettsville; Stone bridges--Tennessee--Goodlettsville; Architectural elements--Tennessee--Goodlettsville	WSC 41	edit Old Stone Bridge in Goodlettsville, Tennessee, n.d.	Charles Warterfield Architectural Collection	An undated photograph the Old Stone Bridge over Mansker's Creek in Goodlettsville, Tennessee. This historic bridge is one of Tennessee's oldest bridges made of stone with masonry arch construction. This dual arch bridge was built circa 1837-1849 as part of the stagecoach road that connected Nashville to Louisville. The Old Stone Bridge is located about 12 miles north of Nashville. Photographed by the Nashville architect Charles Wesley Warterfield, Jr. (1926-1998); the architectural photograph forms part of the Charles Warterfield Architectural Collection. Gift of: The Estate of Charles W. Warterfield, Jr., F. A. I. A. The original is a 35 mm color transparency (2 x 2 in. slide mount).	St Pl Ar ph
Goodlettsville (Tenn.)--Buildings, structures, etc.; Goodlettsville	<a href="http://nashville.contentdm.oclc.org/cdm/ref/collection/nr/id/4">http://nashville.contentdm.oclc.org/cdm/ref/collection/nr/id/4</a>				

# Hydra & F3

- Overlapping Datastreams to PCDM Objects
- RELS-EXT != PCDM Relationships
- Keep that MODS or XML Metadata in Hyrax As Needed In Your Transition
- Fedora Migrate ... ?  
<https://github.com/samvera-labs/fedora-migrate>
- FOXML & Fedora File System-based Migration ... ?  
<https://github.com/fcrepo4-exts/migration-utils>
- Consider Administrative Data Placement for Files / Filesets

A photograph of a dirt road winding through a dense forest of tall evergreen trees. The road is paved and curves to the right, disappearing into the distance. The surrounding trees are tall and thin, with many needles and branches. The lighting suggests it might be early morning or late afternoon, with sunlight filtering through the trees.

# You're Not Alone in this Process

# Helpful Resources & Groups

- [Samvera Community Groups](#)
  - [DSpace Migration Group](#)
  - [Documentation Group](#)
  - [Metadata Group](#)
  - [Hyku Metadata Guidance](#)
- Outside of Samvera...
  - Tools:[OpenRefine](#), [Catmandu](#), ...
  - [PCDM Community](#)
  - Our Fedora 4 Peers: [Islandora!](#)
  - [DLF Metadata Assessment Group](#)
- Peer Groups...
  - [Hyrax Renegade Metadata Docs](#)
  - [Shared Cultural Heritage MAPs](#)
  - [DPLA MAP v.4](#)

# Helpful Resources & Groups

- [Samvera Community Groups](#)
  - [DSpace Migration Group](#)
  - [Documentation Group](#)
  - [Metadata Group](#)
  - [Hyku Metadata Guidance](#)
- Outside of Samvera...
  - Tools:[OpenRefine](#), [Catmandu](#), ...
  - [PCDM Community](#)
  - Our Fedora 4 Peers: [Islandora!](#)
  - [DLF Metadata Assessment Group](#)
- Peer Groups...
  - [Hyrax Renegade Metadata Docs](#)
  - [Shared Cultural Heritage MAPs](#)
  - [DPLA MAP v.4](#)

**And Please  
Remember to  
Share Back Your  
Experiences &  
Lessons as Well!**



**Thank you!**

**<http://bit.ly/HyraxDataMigration>**