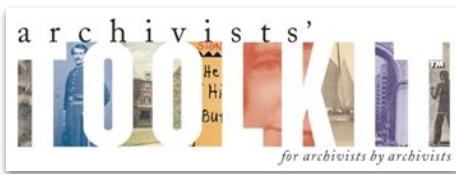COLUMBIA UNIVERSITY LIBRARIES

# CUL ArchivesSpace Migration Project

**David Hodges and Kevin Schlottmann**

**Project phases:**

1. Migrate collection-level data (~8,000 collection and accession records), and integrate with Voyager ILS (2017–18)

2. Remediate and migrate ~1,400 EAD finding aids (2019)

3. Implement public-facing interface for finding aids, and integrate with other discovery services (2019–20)

```
100   1b     $aProvenance, William Fonds,$d1897-1938.
245   00     $kPapers,$f[ca. 1917-1937].
```

**Data and workflow complexities:**

- Several archival repositories (RBML/Univ. Archives, Avery, Burke, Starr) with divergent tools/workflows

- ~4K MARC records

- ~4K accessions in Archivists' Toolkit

- ~1,400 legacy EAD finding aids

- Bespoke finding aids in PDF or HTML

- Linked spreadsheets and other binary formats

```
● ldpd_4079753_ead.xml ×

processing-instruction
     1    <?oxygen NVDLSchema="http://eadrepo.cul.columbia.edu:8080/exist/rest/db/ead/nnc-rb_staging/schema/ead_cul.nvdl" type="xml"?>
     2 ▽  <ead xmlns="urn:isbn:1-931666-22-9" xmlns:xlink="http://www.w3.org/1999/xlink" audience="external">
     3 ▽      <eadheader findaidstatus="draft" countryencoding="iso3166-1" dateencoding="iso8601" langencoding="iso639-2b" relatedencoding="DC" repositoryencoding="iso15511" scriptencoding="iso15924">
     4            <eadid countrycode="US" encodinganalog="Identifier" publicid="-//us:nnc-rb//TEXT us::nnc-rb::ldpd_4079753_ead//EN" mainagencycode="nnc-rb">ldpd_4079753_ead.xml</eadid>
     5 ▽          <filedesc>
     6 ▽              <titlestmt>
     7                    <titleproper encodinganalog="Title">Guide to the Carnegie Corporation of New York records,</titleproper>
     8                    <author encodinganalog="Contributor">Columbia University Libraries. Rare Book and Manuscript Library</author>
     9                </titlestmt>
    10 ▽              <publicationstmt>
    11                    <publisher encodinganalog="Publisher">Rare Book and Manuscript Library<lb/>Columbia University Libraries</publisher>
```

**Motivations:**

- Single-source of truth for description
- Improved creation and editing of archival description
- Unify practices and workflows across CU libraries
- Modernize infrastructure
- Improve discoverability and UX

**Phase 2 EAD Migration**

- Combine container structure (<dsc>) from finding aids with collection-level data already in AS and MARC records.

- Some collection-level description was better in legacy finding aids, but varied from record to record.

- How to determine the source of truth for conflicting data? Review process could be extremely time consuming.

**EAD_Process (Python + Google Sheets)**

- Python + Google Sheets + XSLT: extract data, flag elements to merge, process all files.

- Two data sources: legacy finding aid and EAD exported from AS processed in pairs.

- Defined list of XPATH queries for each.

- Data is extracted from both data sets and sent to a Google Sheet, which collates into side-by-side view for each of targeted elements.

```python
# Dict of elements and their xpath. These are used for legacy EADS.
legacyQs = {
    "bibid": "ead:archdesc/ead:did/ead:unitid[@type='clio'][1]/text()",
    "repo": "ead:archdesc/ead:did/ead:unitid[1]/@repositorycode",
    "title": "ead:archdesc/ead:did/ead:unittitle[1]/text()",
    "status" : "ead:eadheader/@findaidstatus",
    "revisiondesc": "ead:eadheader/ead:revisiondesc",
    "altformavail": "ead:archdesc/ead:altformavail",
    "accruals": "ead:archdesc/ead:accruals",
    "accessrestrict": "ead:archdesc/ead:accessrestrict",
    "userestrict": "ead:archdesc/ead:userestrict",
    "acqinfo": "ead:archdesc/ead:acqinfo",
    "arrangement": "ead:archdesc/ead:arrangement",
    "bibliography": "ead:archdesc/ead:bibliography",
    "bioghist": "ead:archdesc/ead:bioghist",
    "scopecontent": "ead:archdesc/ead:scopecontent",
    "controlaccess": "ead:archdesc/ead:controlaccess",
    "custodhist": "ead:archdesc/ead:custodhist",
    "separatedmaterial": "ead:archdesc/ead:separatedmaterial",
    "otherfindaid": "ead:archdesc/ead:otherfindaid",
    "relatedmaterial": "ead:archdesc/ead:relatedmaterial",
    "abstract": "ead:archdesc/ead:did/ead:abstract",
    "physloc": "ead:archdesc/ead:did/ead:physloc",
    "processinfo": "ead:archdesc/ead:processinfo",
    "unitid": "ead:archdesc/ead:did/ead:unitid",
    "prefercite": "ead:archdesc/ead:prefercite"

}
```

```python
def ead_report(the_sheet,the_range,the_files,the_qs):

    # Default namespace for CUL EADs
    ns = {"ead": "urn:isbn:1-931666-22-9"}

    the_heads = list(the_qs.keys())
    the_xpaths = list(the_qs.values())
    the_data = [the_heads,the_xpaths]


    gs.sheetClear(the_sheet,the_range)

    # this will be the result dict.
    theElements = { }

    for a_file in the_files:

        print('Processing file: ' + a_file)
```

**EAD comparison and review workflow (Google Sheets)**

- Each element is presented side by side for comparison.

- Archivists review and flag which is the "correct" one.

- Some heuristics help automate decision, e.g., if text size differs by +x chars, then migrate the longer.

**EAD_Merge (Python + Google Sheets + XSLT)**

- After archivists have completed review, a Python script reads Google Sheet migration grid rows as lists ([['4079432', 'bioghist', 'scopecontent', 'abstract'],...]).

- Legacy EAD file is sent to XSLT with parameters of which additional elements to migrate (other than <dsc>).

- XSLT merges two EAD trees, incorporating <dsc> plus selected additional elements per the parameters.

- Another XSLT stylesheet in pipeline performs other global cleanup functions.

- Output is validated and QC'd before importing into AS.

```python
def get_migration_grid(theSheet,theRange):

    the_data = []

    x = gs.getSheetData(theSheet, theRange)

    the_values = x["values"]
    the_heads = the_values[0]

    for a_row in the_values:
        my_bibid = a_row[0]
```

```xml
<!-- Capture the <dsc> section of the source tree for later use. -->
<xsl:variable name="the_dsc">
    <xsl:copy-of select="//dsc"/>
</xsl:variable>
```

```xml
<!--Special collection-level stuff to optionally migrate replacing existing elements:  -->

<xsl:template
    match="revisiondesc[not(ancestor::dsc)] |
    bioghist[not(ancestor::dsc)] |
    scopecontent[not(ancestor::dsc)] |
    relatedmaterial[not(ancestor::dsc)] |
    prefercite[not(ancestor::dsc)] |
    custodhist[not(ancestor::dsc)] |
    acqinfo[not(ancestor::dsc)] |
    processinfo[not(ancestor::dsc)] |
    accessrestrict[not(ancestor::dsc)] |
    userestrict[not(ancestor::dsc)] |
    abstract[not(ancestor::dsc)]
    "
    mode="asead">
    <xsl:call-template name="replaceElement">
        <xsl:with-param name="theElement">
            <xsl:value-of select="local-name(.)"/>
        </xsl:with-param>
    </xsl:call-template>
</xsl:template>
```

# *Thank you!*

⇔

## Code repo

https://github.com/cul/rbml-archivesspace

## Contact info

• @archivistkevin // kws2126@columbia.edu

• dwh2128@columbia.edu