# Cleaning up legacy metadata for ETDs:
## Strategies, tools and a look into the future

Xiping Liu, Albert Duran, Anne Washington

UNIVERSITY of HOUSTON | LIBRARIES
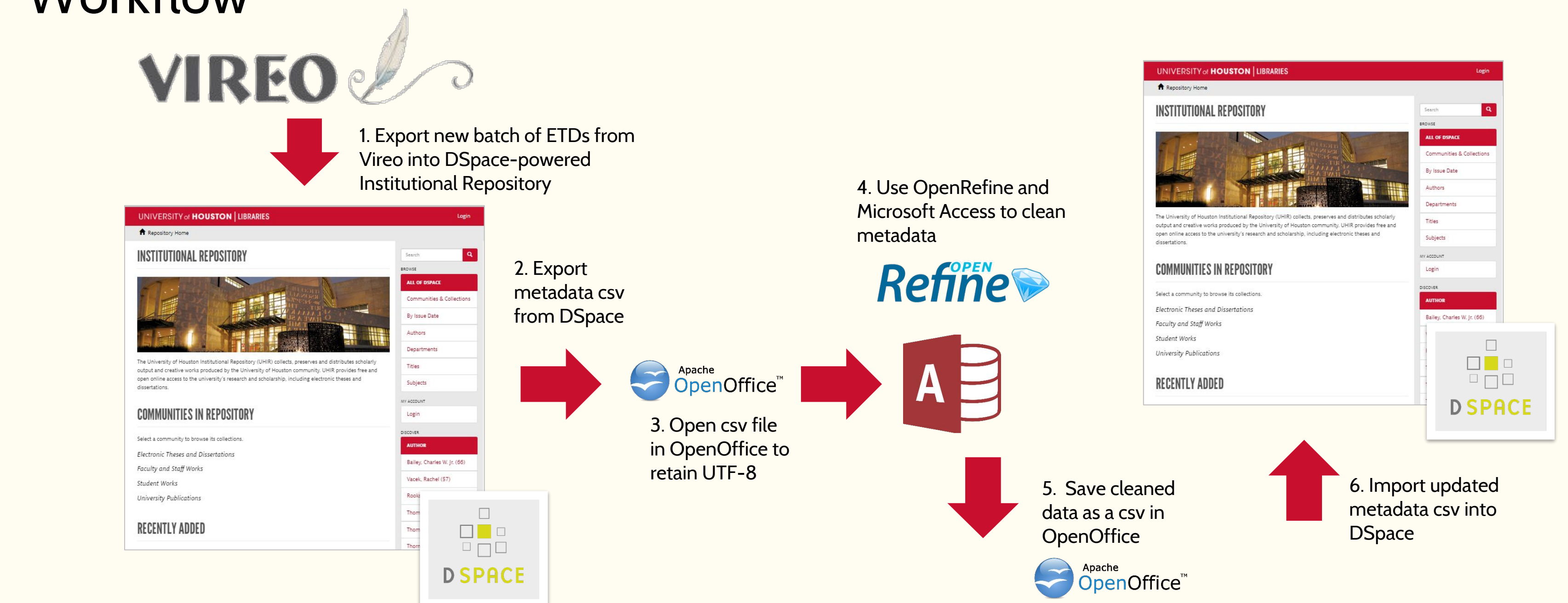
## Background

Since July 2015, the University of Houston (UH) Libraries Metadata and Digitization Services (MDS) department has collaborated with the Digital Repository Services (DRS) department to improve the quality of legacy Electronic Thesis and Dissertation (ETD) metadata in the UH Institutional Repository. The team developed workflows and strategies that improved the ETD metadata quality, strengthened the communication and collaboration between the DRS department and the MDS Metadata Unit, and provided insight into opportunities for future development.

## ETD Metadata Upgrade Goals

- Align metadata with Texas Digital Library (TDL) standards*
- Correct department and discipline names
- Standardize advisor and committee member names
- Normalize date fields
- Remove broken URLS
- Develop sustainable workflow for bi-annual additions to DSpace

*http://hdl.handle.net/2249.1/68437

## Workflow



1. Export new batch of ETDs from Vireo into DSpace-powered Institutional Repository

2. Export metadata csv from DSpace

3. Open csv file in OpenOffice to retain UTF-8

4. Use OpenRefine and Microsoft Access to clean metadata

5. Save cleaned data as a csv in OpenOffice

6. Import updated metadata csv into DSpace

## Communication and Documentation



Basecamp facilitates interdepartmental collaboration between the MDS Metadata Unit and Digital Repository Services. Basecamp can issue to-dos, track progress, save discussions, and maintain a record of the minutes from weekly team meetings.

MDS uses a departmental wiki, powered by PmWiki, to document workflow processes and archive project information such as a List of colleges, departments and degree programs that submit ETDs to the UH IR. PM Wiki is free and Open Source Software.

## Remediation



DSpace metadata exports produced repeated fields that needed to be collapsed together. A Microsoft Access query was used to consolidate values from the duplicate columns into a single column. This consolidation is necessary to complete additional remediation tasks and import the new metadata successfully into DSpace.

OpenRefine was used to standardize advisor and committee member names. All names were collapsed into a single column using the transpose function. Then, the facet and cluster functions were used to choose and populate the preferred name throughout the dataset.

## Re-import



Once all of the metadata clean-up tasks were complete, the updated metadata csv files were imported into DSpace and records were updated.

## Where Are We Now?

- ✓ 941 corrected records
- ✓ Workflow for new batches
- ✓ Corrections to Vireo



Map by Karsten Barnett from the Noun Project

## Next Steps

- Add names to local thesaurus
- Workflow efficiencies
- Data sharing