# The ArchivesDirect pilot: Road-testing Archivematica hosting in DuraCloud

Amy Rushing (UTSA), Julianna Barrera-Gomez (UTSA) and Courtney C. Mumma (Artefactual Systems)
TCDL, Austin, TX, April 27, 2015, Session 4A

# Overview

- The pilot
- Administrative experience - UTSA partner
- Operational experience
- Q&A and discussion

# What were the pilot's goals?

- standards-based digital preservation packages in secure long-term storage
  - Archivematica + DuraCloud
  - hosted, web-based solution
  - customizable to match your needs/workflows/content
  - no vendor/content lock-in
- identify service needs and scope

# Who were the pilot volunteers?

Berea College

Huntington Library

Illinois-Wesleyan University

Kansas State University

North Carolina Dept of Cultural Resources

Pepperdine University

Phillips Academy Andover

University of Texas at San Antonio

University of Washington

# What was the pilot structure?

- assess workflows for distinct content
  - preservation and/or access copies? service masters? manual normalization?
  - transcription, forensic analysis, directory structure for OO?
  - format identification tool, error handling, compression and extraction preferences?

- assess content
  - born-digital, digitized, MD and submission documentation
  - size and scale

# Structure, cont'd

# Structure, cont'd

# Structure, cont'd

- intensive training and content analysis guidance
- pilot wiki for sharing workflows and results
- pilot-support and pilot-discussion lists
- agile development processes with rapid analysis/ upgrade/testing cycles

# What were the pilot outcomes?

- ArchivesDirect public offering and price structure
- toolkits and documentation
- DuraCloud integration and scalability enhancements to Archivematica
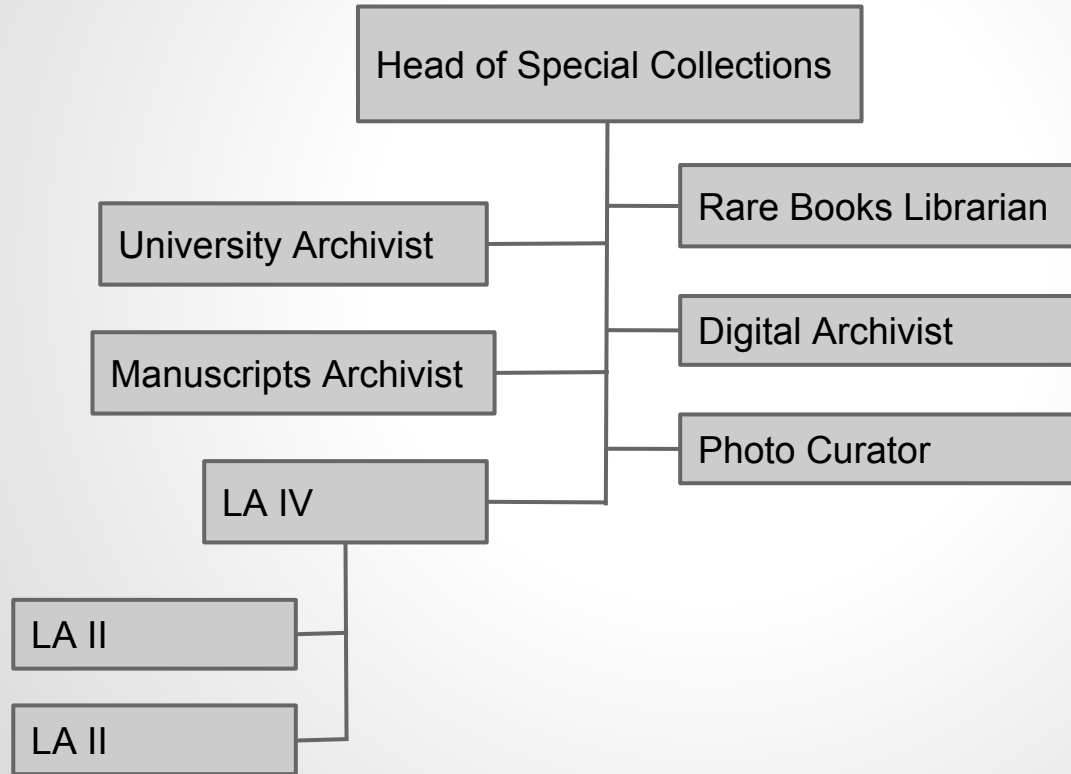- training and support strategy



archivesdirect.org

# Special Collections

# IT @ UTSA

| Estimated network drive content sizes by type: | | Total Size |
|---|---|---|
| | | |
| **At-risk digitized content:** | | |
| Images\Marquise - JPEGS | 7.46 GB | |
| thanatus\NEDCC | 5.41 TB | |
| thanatus\NEDCC_Batch2 | 1.86 TB | |
| thanatus\NEDCCStagingArea | 602 GB | |
| thanatus\Portal of Texas History newspapers | 789 GB | |
| thanatus\PTH Test | 179 GB | **8.85 TB** |
| | | |
| **At-risk born-digital content:** | | |
| archives_data\Archives Archive | 5.04 GB | |
| archives_data\COLLECTIONS | 99.6 GB | |
| archives_data\Electronic Records | 607 GB | |
| archives_data\SARA | 36 MB | |
| library_data\UA 16.01.01 | 508 GB | **1.26 TB** |
| | | |
| **General digitized content:** | | |
| UTFILE\Audio | 124 GB | |
| UTFILE\Media Content on Helix Server | 101 GB | |
| UTFILE\Video | 558 GB | |
| archives_data\Images (-7.46 GB from Marquise) | 994 GB | |
| library_data\master | 1.91 TB | |
| library_data\_I Drive - Common\PhotoScan | 10 GB | **3.70 TB** |
| | | **13.81 TB** |
| | | |
| **Removable media inventories:** | | Total Size |
| MS collections | 4 TB | |
| UA collections | 454 GB | **4.45 TB** |
| | | |
| **Estimated Grand Total:** | | **18.26 TB** |

# Jump in Initiative

# Jump In Manuscripts Inventory



Quantity of Media by Format

# Jump In University Archives Inventory



Quantity of Media by Format

CD    DVD    3.5" disk    5.25" disk    Zip disk

388, 61%

171, 27%

66, 10%

12, 2%

4, 1%

8, 1%

# Pilot Content overview

*San Antonio Light* digitized negatives

Tiffs and jpegs, plus metadata in .csv

Born Digital material from Dr. Ellen Riojas Clark

Facebook export, Outlook email

Born Digital material from Dean George Perry

.msg email files with attachments, Word files

Born Digital material from Dean Dan Gelo

Word and .pdf

*The Marquise* digitized newspapers

jpegs

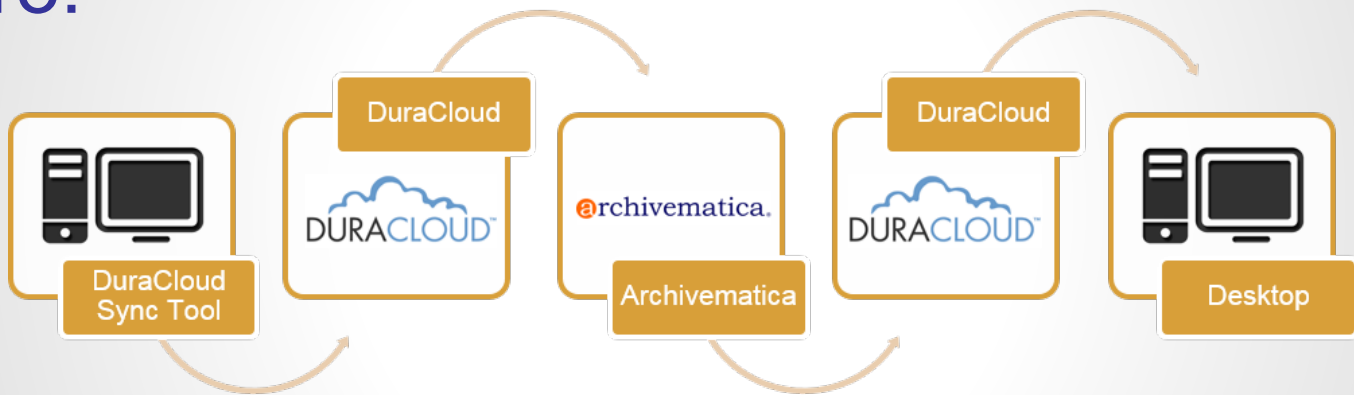Born Digital architecture files from Killis Almond & Associates

variety of CAD file types

# Processing POV

# Workflow

Macro:

# Workflow con't.

Micro:
- Set up
- Assemble testers
- Coordinate tests
- Generate transfers, SIPs and AIPs
- Track results
- Track issues

# Test log



Archivematica Test Log

| | Testing Dates | Name | Test Directory | Notes |
|---|---|---|---|---|
| 1 | **Testing Dates** | **Name** | **Test Directory** | **Notes** |
| 2 | Dec. 11-12 | Julianna | MexicanCookbooks\txsau_tx714-m47-1930z | During Ingest noticed same directory appears to have been run through on 2014-10-30 |
| 3 | Dec. 15-16 | Amy | NEDCC\tif 1 | |
| 4 | | Amy | KillisAlmond\MS348 Hardrive\UTSA Digital Back Ups\Digital Project Files from Archive\1894 GRAND OPERA HOUSE\OLD SHEETS AND SCANS | |
| 5 | | Amy | EllenClark/ClarkOutlook | |
| 6 | Dec. 17-18 | | | |
| 7 | Jan. 5-7 | Julianna | Perry\Correspondence\2006\L-Z | Changed from FIDO to file ext in both the Transfers and Ingest file ID configuration |
| 8 | Jan. 8-9 | Julianna\Amy | Perry\Correspondence\2006\G-K | Will change from FIDO in the Dashboard to aid in format IDing for this transfer of .msg files, then switch back, per Justin Simpson's instructions |
| 9 | Jan. 12-13 | Amy | Killis Almond\....06-07 - Ft. Sam | Amy was unable to do testing this week |
| 10 | Jan. 14-15 | | | |
| 11 | 2/12/2015 | Julianna | KillisAlmond | Complete Hard drive |
| 12 | 2/16/2015 | Julianna | Marquise\Marquise | |
| 13 | | | Marquise\Marquise paste-ups | |
| 14 | 2/5/2015 | Julianna | Gelo\Acc.2011-005 | |
| 15 | 2/5/2015 | Julianna | EllenClark\ClarkFacebookExport | |
| 16 | 2/11/2015 | Julianna | CDMUpload002FebTry2 | ALTERED the Administration setting so that normalize is manual only. This is set up with the manualNormalization directory and also includes a metadata.csv file |
| 17 | 2/12/2015 | Julianna | CDMUpload002FebTry3 | Admin settings: do not normalize. Will try this, since access & pres are already set up. |
| 18 | 2/17/2015 | Julianna | MexicanCookbooks\txsau_tx714-m47-1930z | Running this again as a test for the Ingest: Transcribe SIP microservice OCR txt file production |

# Test log: issues



**Archivematica Test Log**

| SpecColl Staff Name | Transfer Date | Transfer Name | Transfer # Objects | Transfer Size | Transfer Time Started | Ingest Date | Ingest Time Started | Error Reporting : Micro-service where error occurred; what job it was; paste in exact wording |
|---|---|---|---|---|---|---|---|---|
| Julianna | 12/11/2014 | MXCookbookstx714m47 | 366 jpegs | 489 MB | 4:43 PM | 12/11/2014 | 17:01 | Micro-service: Store AIP Job: Store the AIP Failed |
| Amy Ru: | 12/15/2014 | NEDCC small batch 1 | 5 tiffs | 1.03 GB | 10:39 AM | 12/15/2014 | 10:45 | Micro-service: Transcribe SIP Contents Job: Transcribe Failed |
| Amy Ru: | 12/16/2014 | Killis Almond-OperaHouse-OldSheetsandScans | 71 files of various formats (DWG, txt, pdf, CAL, bmp) | 81 MB | 13:46 | 12/16/2014 | 13:53 | |
| Amy Ru: | 12/16/2014 | ClarkOutlook-SecondTry | 1 .pst file | 59 MB | 2:45 PM | 12/16/2014 | 2:45 PM | Micro-service: Normalize Not normalizing ClarkEmailAcc2014-20.pst - No rule or default rule found to normalize for preservation |
| Julianna | 1/5/2015 | PerryemailLtoZ | 5510 .msg files | 515 MB | 2:07 PM | 1/5/2015 | 17:27 | Micro-service: Identify file format Job: Identify file format (transfer continuing) ; Micro-service: Characterize and extract metadata Job: Characterize and extract metadata Failed (transfer continuing) ; INGEST: Micro-service: Normalize Job: Identify file format; |
| Julianna | 1/8/2015 | PerryemailGtoK | CANCELLED | CANCELLED | CANCELLED | CANCELLED | CANCELLED | Rejected transfer--accidentally picked wrong folder from Perry folder (picked L-Z) |
| Julianna | 1/8/2015 | PerryemailGtoKretry | 4199 .msg files | | 4:31 PM | 1/9/2015 | 12:01 | No errors!  Turned format id method back to FIDO in dashboard. |
| Julianna | | KillisAlmondCompleteHD | apprx.1,800 files | 954 MB | 12:31 PM | 2/12/2015 | 5:10 PM | Transfer: Micro-service: Identify File Format, Job: Identify File format; Micro: Extract packages: Job: Identify file format; Micro: Chararacterize and extract metadata: Job: Characterize and extract metadata;; Ingest: Micro: Normalize, Job: Identify file format; Micro: Transcribe SIP contents, Job: Transcribe |
| Julianna | 2/16/2015 | MarquiseFeb | 1,587 jpegs | 3.03 GB | 8:33 AM | 2/16/2015 | 11:03 AM | {Pipeline stuck on Ingest:Transcribe SIP contents--Emailed pilotsupport, will take several hours} FAILED Ingest: Microservice: Prepare AIP, Job: Prepare AIP (failed, email report sent automatically) |

# Support methods



- ## ArchivesDirect:
  - DuraSpace group wiki
  - Pilot support email list
  - Consultations


- ## Archivematica:
  - User manual 1.2 wiki
  - Archivematica Google group

# Testing experience

## Assembling transfers

- Metadata set-up
- Directory structure
- DuraCloud Sync tool

# Archivematica testing

- Transfers
- Micro-services
- Specific tools
- Processing methods

# Outcome

Errors & issues

# Outcome

## Errors & issues

- File-level errors

# Outcome

## Errors & issues

- File-level errors
- Normalization errors

# Outcome

## Errors & issues

- File-level errors
- Normalization errors
- Metadata errors

# Outcome

## Errors & issues

- File-level errors
- Normalization errors
- Metadata errors
- Fatal errors



Archivematica Fail Report for Transfer: bag-9827956f-4af1-4579-b62c-bdf983d52238

Trash   x

ArchivematicaSystem@archivematica.org via artefactual.com        Apr 4    Reply

to courtney

This message has been deleted. Restore message

| unitType | Total time processing | total file size | number of files | average file size KB | average file size MB |
|---|---|---|---|---|---|
| Transfer | 0:20:24 | 12403961 | 34 | 364.82238235 | 0.36482238 |

| Type | Status | Started | Duration |
|---|---|---|---|
| Verify bag, and restructure for compliance | Failed | 2013-04-04 23:12:57 | |
| Assign checksums and file sizes to objects | Completed successfully | 2013-04-04 23:12:48 | |
| Assign file UUIDs to objects | Completed successfully | 2013-04-04 23:12:36 | |
| Rename with transfer UUID | Completed successfully | 2013-04-04 23:12:35 | |
| Move to processing directory | Completed successfully | 2013-04-04 23:12:33 | |
| Set file permissions | Completed successfully | 2013-04-04 23:12:31 | |
| Approve bagit transfer | Completed successfully | 2013-04-04 23:12:20 | |

# **Outcome**

## Errors & issues

- File-level errors
- Normalization errors
- Metadata errors
- Fatal errors
- General pain points

Connected ●

# Final thoughts/Next steps



- Piloting alternative ArchivesDirect workflow (not yet in production)
- Positive pilot outcome will result in alternative hosting option for ArchivesDirect users

# helpful links

archivesdirect.org

lib.utsa.edu/special-collections

Archivematica Birds-of-a-Feather Session (with boxed lunch) - PCL Map Room, Tues 12-130

archivematica.org - duracloud.org - artefactual.com - duraspace.org