# Taming Cerberus:
## Balancing Legacy Digitization and Emerging Needs

**Carnegie Mellon University** | March 2022

Julia Corrin, Ann Marie Mesco

# Collections by the Numbers

**26** years

**2.2** million pages

**37+** terabytes

# Digitization is more than just scanning…

- Manage legacy collections to ensure ongoing accessibility

- Adapt to evolving standards and technologies for new collections

- Scanning as a skill set



*Hans Sebald Beham, Public domain, via Wikimedia Commons*

# Access is more than just a platform…

- User expectations & knowledge
  - Different access needs for different audiences
- Accessibility
  - WCAG, ADA
  - Annotated PDFs
- Metadata
  - Search and retrieval is limited to what a computer can understand

# Whither Metadata?

- Metadata is under-discussed when it comes to digitization

- Metadata is frequently the most resource intensive aspect of digitization

- Metadata can be the source of most long-term technical debt issues for digital projects

# Case Study 1: Sen. H. John Heinz III Papers

**Goal**: Replicate the experience of using the physical collection in a digital space.

- Received by the archives in 1991, digitized in 1994
- Over $1.5 million in project funding
- Over 300,000 documents / 390 cubic feet

# 1994 Challenges - The Wild West

- Lack of best practice
  - No established scanning guidelines for this type of scanning
    - FADGI - 2007
  - Limited understanding of metadata
    - DublinCore - 1995; EAD 1.0 - 1998
- Technology Limitations
  - Scanners were expensive and limited in what they could do
    - 400 DPI, black and white TIFFs
  - No widely adopted content management systems

**HelioScan v1.1**

| HELIOS SUBGROUP | Legislative Records |
| SERIES | Committee Files -- 1977-1991 |
| SUBSERIES | Aging Committee -- 1980-1991 |
| SUB-SUBSERIES | Professional Staff Member Files -- 1982-1991 [1987-1991] |
| SUB-SUB-SUBSERIES | Fiegener, Janice -- 1982-1991 [1989-1990] |

BOX [NEW] [EDIT] 123
FOLDER [NEW] [EDIT] 1: Social Security Administration -- Attorney Fee Bill -- 1991
BUNDLE ☒ [NEW] 0001
DOCUMENT [NEW] 0001: Memoranda -- 1989
PAGE ☒ display [<<] 000001 [>>]

REVISION [RESCAN (ADF)] [RESCAN (FB)] 01

SIZE: ○ Letter ○ Legal
SIDES: ● Single Sided ○ Double Sided
QUALITY: ○ Faint ● Normal ○ Dark
ORIENTATION: ● Portrait ○ Landscape

[Scan from Feeder] [Scan from Flatbed] [Help]

# 2022 Challenges

- TIFFs can no longer be read by modern TIFF libraries

    - Required to use modern viewers like Mirador

- Periodic re-OCR to leverage technology improvements

    - Re-OCRing 700,000 pages is resource intensive

- Long term file management

    - Master vs. derivative

    - Corruption overtime

# More 2022 Challenges

- Metadata remediation
  - Increased "noise" over time
- Structural assumptions
  - The way users interact with the collection has evolved
  - 1994 assumptions about users are baked into the collection
- Persistency
  - DOIs, PURLs need to be maintained overtime
  - Limits changes that can be made

## Metadata

**Browse**
All Documents > H. John Heinz III > Legislative Directors' Files -- 1977-1991 (1979-1981, 1987-1990) > >browse3_ss:"Civil Rights">Civil Rights > >browse3_ss:"Civil Rights">>browse5_ss:"Civil Rights Act of 1990 -- JH Working Files">Civil Rights Act of 1990 -- JH Working Files

**Title**
-- 1991 (bundled) (Civil Rights -- Civil Rights Act of 1990 -- JH Working Files -- 1990)

**Collection**
H. John Heinz III

**Series**
Legislative Directors' Files -- 1977-1991 (1979-1981, 1987-1990)

**Archival Topic**
Civil Rights

**Folder Title**
Civil Rights Act of 1990 -- JH Working Files

**Identifier**
\Heinz\box00326\fld00006\bdl0007\doc0002\Heinz_box00326_fld00006_bdl0007_doc0002.p

**Rights**
Legislative Records -- 1970-1991 (1977-1991)

**Type**
pdf

**Thumbnail**
\Heinz\box00326\fld00006\bdl0007\doc0002\THUMBNAIL\Heinz_box00326_fld00006_bdl00

**Document ID**
734887
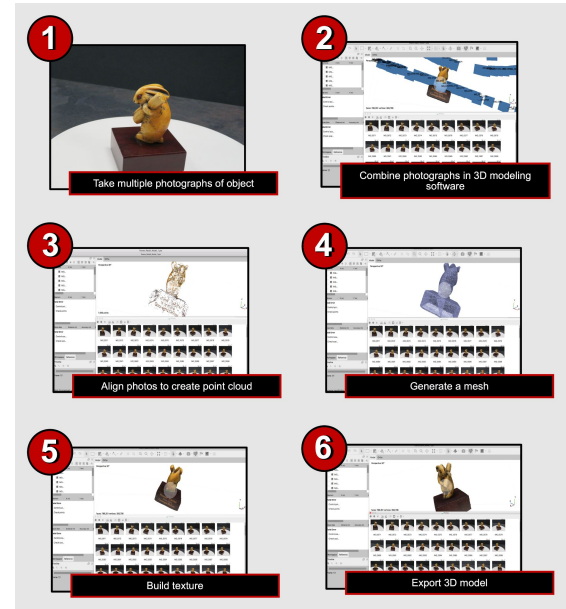
# Case Study 2: New Formats

**Goal:** Capture "non-traditional" theses, dissertations, and research outputs

- Equitable preservation of scholarly output, regardless of discipline
- Expand scanning capabilities without expensive purchases
- Types of objects: Architectural models, design projects

# Method 1: Photogrammetry

- Create a 3D model from a series of 2D photographs

- Lower cost entry point than dedicated 3D scanners

- Preservation advantages

- Leverages technology and skills already available in many digitization labs

**Carnegie Mellon University**

# Photogrammetry Starter Kit

- Decent digital camera
  - Cannon DSLR EOS 80D
- Computer
  - Digital Storm Velox
  - Gaming computer, VR Computer
- Software
  - Agisoft Metashape
  - Meshmixer
- Turntable

Perspective 30°

Snap: Axis, 30



faces: 565,677 vertices: 283,349

Photos

IMG_3461    IMG_3462    IMG_3463    IMG_3464    IMG_3465    IMG_3466    IMG_3467    IMG_3468    IMG_3469    IMG_3470    IMG_3471    IMG_3472    IMG_3473    IMG_3474
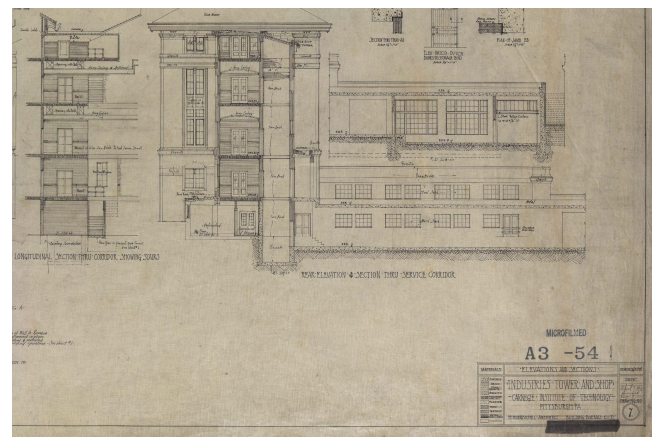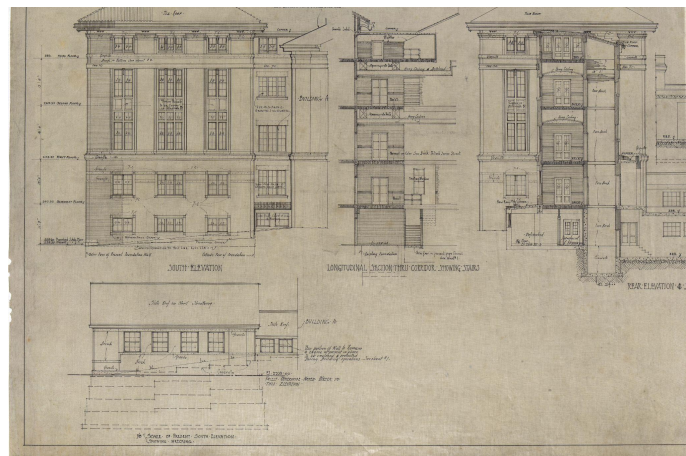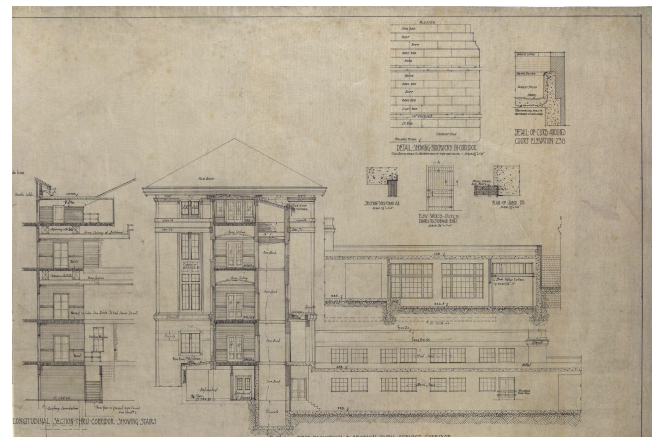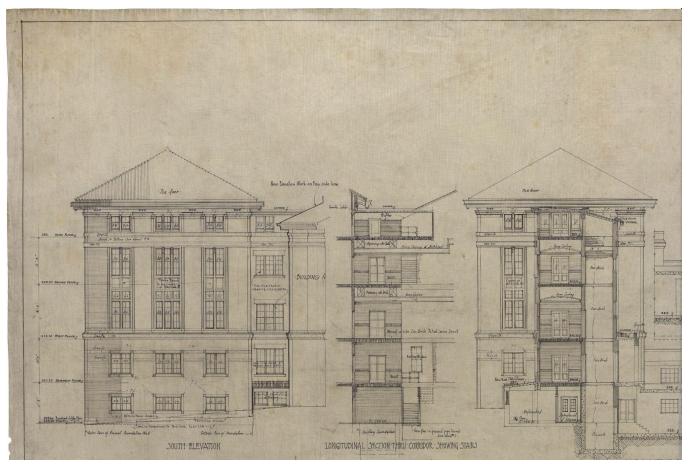
Photos  Console  Jobs
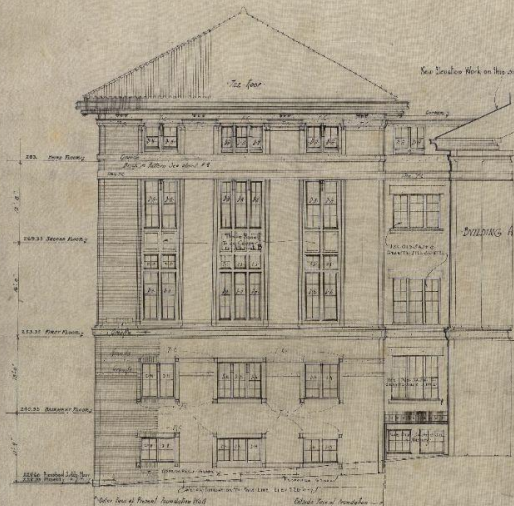
# Method 2: Image Stitching

- Capture large format works that will not fit on a traditional scanner
- Digitize fragile large format works without damaging the original
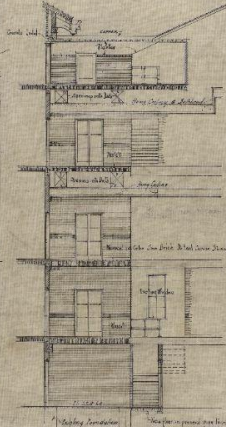- Utilize skills of experienced scanner operators

# Large Format Starter Kit

- Existing overhead scanner

  - CopiBook scanner; 23x16 in scan area

  - Large documents require 4-6 captures
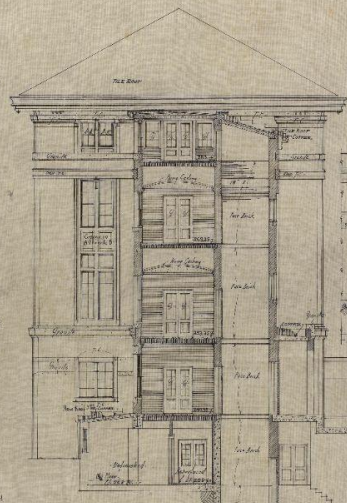
- Adobe Photoshop

- People!

A3 -54

SOUTH ELEVATION

LONGITUDINAL SECTION THRU CORRIDOR SHOWING STAIRS

REAR ELEVATION & SECTION THRU SERVICE CORRIDOR

DETAIL OF CURB AROUND COURT ELEVATION 2750

DETAIL SHOWING BRICKWORK IN CORRIDOR

DETAIL OF CURB AROUND
COURT ELEVATION 238

BUILDING A

SOUTH ELEVATION

LONGITUDINAL SECTION THRU CORRIDOR SHOWING STAIRS

REAR ELEVATION & SECTION THRU SERVICE CORRIDOR

BUILDING A

⅛ SCALE OF PRESENT SOUTH ELEVATION
SHOWING WRECKING

A3 -54

ELEVATIONS AND SECTIONS
INDUSTRIES TOWER AND SHOP
CARNEGIE INSTITUTE OF TECHNOLOGY
PITTSBURGH PA

7

# Conclusions

- Don't start a digitization project without a metadata plan in place
- Always consider long term needs - every new digitization project reduces your bandwidth overtime
- Invest in people, not just technology

# Acknowledgements

Ann Marie Mesco, Digitization Projects Manager

Joe Mesco, Scanner Operator

Jon McIntire, Scanner Operator

Lina Crowe, Metadata Specialist