

A Novel Workflow for Large Scale Thesis Digitization

TCDL 2016



The rising STAR of Texas

A Novel Workflow for Large Scale Thesis Digitization

**Texas State University
Albert B. Alkek Library**

Jeremy Moore - Digital Media Specialist

Jason Long - Programmer Analyst

Todd Peters - Head, Digital & Web Services

Why Digitize?

- ❖ Faculty interest
- ❖ Space issues

Library infrastructure renovation- Archives losing ~500 linear feet

ETDs required 201 linear feet of shelf space

3 copies -preservation microfilm copy, Archives hard copy, and a hard copy in the stacks.

How Many?

- ❖ Starting with theses that have a silver halide microfilm preservation copy and second hard copy to debind
- ❖ Approximately 5,000 theses between 1930 - 2010

Since 2005, some theses have been digitized and uploaded into our local Institutional Repository.

More recently, the Graduate College is requiring that all thesis and dissertation be submitted electronically.

How to make available?

- ❖ The copyright and permissions issue is unresolved.
- ❖ No public access at this time.
- ❖ Plan for the future and create digitization workflows that will facilitate batch uploading into the Institutional Repository later.

Novel Process?

1. Prepare physical and digital space
2. Pre-Process physical items
3. Digital Capture
4. Post-Process physical and digital items



Prepare physical and digital space



Algorithmic Project Design

- ❖ Discrete micro-processes apply to physical as well as digital
 - Check
 - Move
 - Check
 - Store
 - Check



Organizing Space

- ❖ Physical and digital space should aid, not hinder the process



- 0.toRescanBW
- 1.toPostProcess-BW
- 1.toPostProcess-Color
- 1a.Problems
- 2.toQC-BW
- 2.toQC-Color
- 3.QCing
- z.Batch2_toAccess
- z.Batch2_toPreservation

Pre-Process physical items



Sacrificing Theses

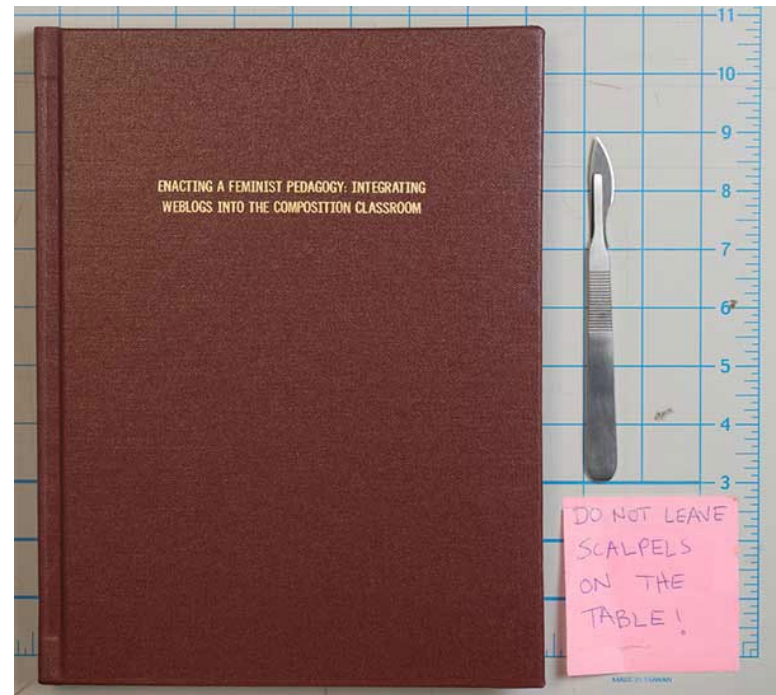
- ❖ **Binding, printing, paper quality changes with each thesis**

- 8.5 x 11 in. in library buckram
- Large margins
- Printed on one side; simplex

- ❖ **Scalpels > Exacto knives**

- Thank you to Jessica Phillips, Head of Preservation at UNT Libraries

- ❖ **Binding, printing, paper quality changes with each thesis**



Chop! Triumph 4315 Stack Cutter

- ❖ Inspect pages before and after chopping
 - Binding, printing, paper quality changes with each thesis
- ❖ Cutting theses to 8.11 in. wide because we had to pick something
 - Binding, printing, paper quality changes with each thesis



Digital Capture



Scanning

- ❖ Fujitsu fi-6670
 - Straight-through feed path
 - Software allows for scanning once and outputting multiple files
 - **Binding, printing, paper quality changes with each thesis**

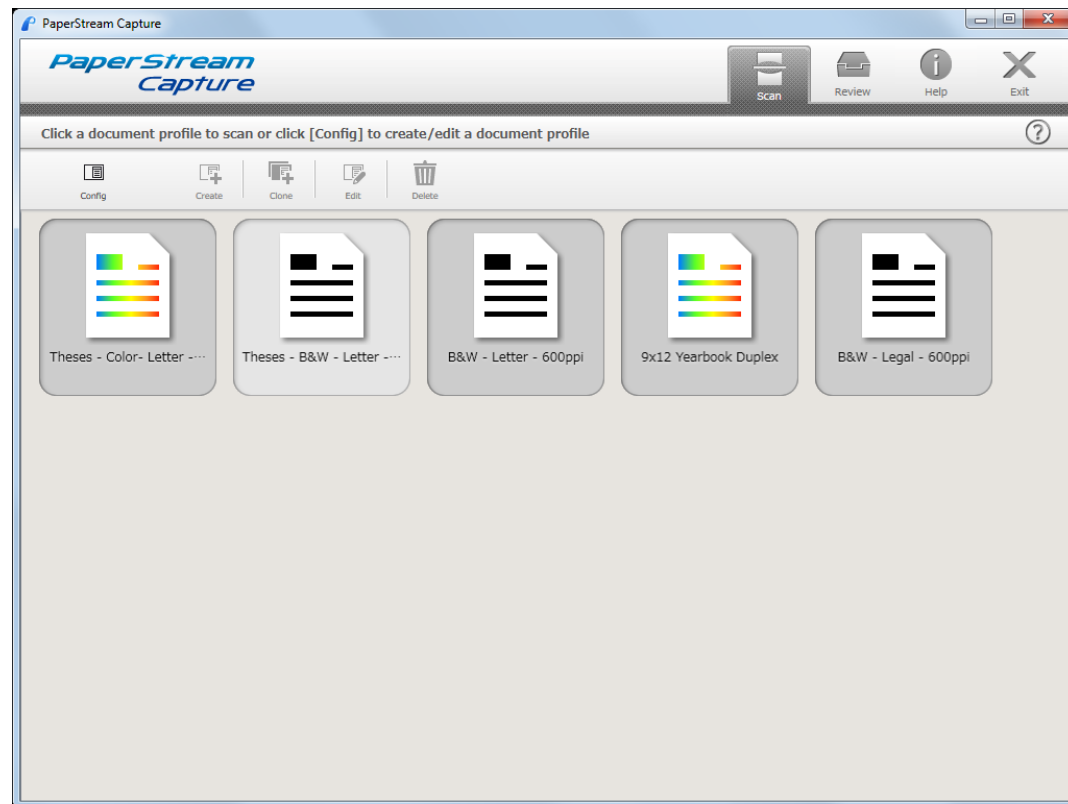
- ❖ MarcEdit Z39.50 client to download Catalog record

- ❖ Name scan folder and MARC record the same:
 - lastName_firstName_year.file
Format



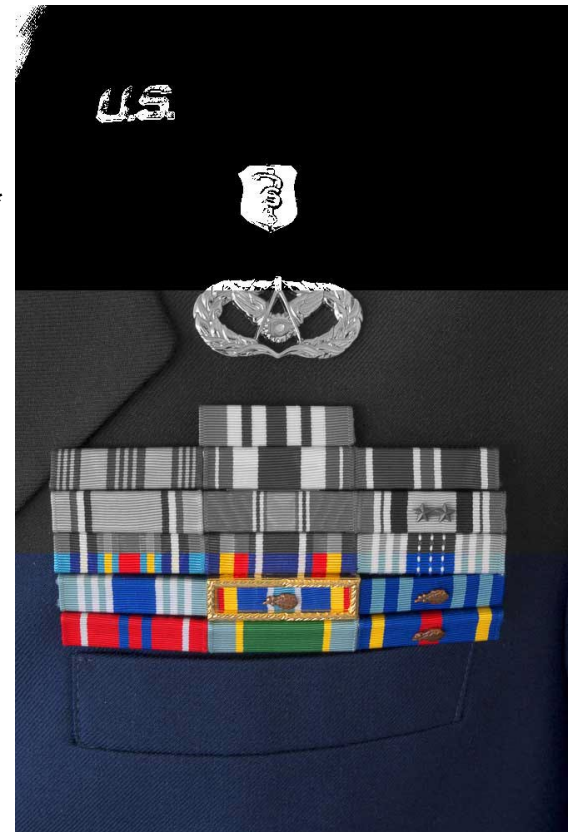
B&W vs Color

- ❖ B&W, or bitonal, is *Only Black and White*
- ❖ *Color* is a catch-all term for *Not Bitonal*



1 bit vs 8 bit vs 24 bit or B&W vs Gray vs Color

- ❖ Black and White (Bitonal) – 1 bit
 - 600 ppi
 - ~4 MB uncompressed TIFF
 - <1 MB group4 TIFF ***SPACE SAVINGS***
- ❖ Grayscale – 8 bit
 - 400 ppi
 - ~15mb uncompressed TIFF
- ❖ Color – 24 bit (3 x 8 bit)
 - 400 ppi
 - ~45 MB TIFF uncompressed



Post-Process physical and digital items



It's about the end product

- ❖ It's about the end product
 - It's about the end product
 - It's about the end product
 - It's about the end product
 - » It's about the end product
 - » It's about the end product
 - » It's about the end product
 - » It's about the end product

» Details matter



Whatever Works (within reason)

- ❖ BASH Shell scripts that will automatically copy, convert, move, name, rename, deskew, resize folders and files
- ❖ Different scripts to process the B&W and Color theses
- ❖ Color pages have a different process than Bitonal
- ❖ Adobe Photoshop and/or ImageMagick to process pages as necessary
- ❖ ABBYY Finereader Hot Folders to batch create and OCR PDFs from Preservation TIFFs

```
#!/bin/bash
```

```
# ProcessBWTheses.sh
```

```
# jeremy.moore@txstate.edu
```

```
# 01-27-2016
```

```
for dir in *; do
```

```
    cd "$dir"
```

```
    echo "$dir"
```

```
    mkdir -pv "$dir"
```

```
    x=1
```

```
    for img in *.tif; do
```

```
        mv "$img" `printf "%s/%s_%04d.tif" \
```

```
            "$dir" "$dir" $x`
```

```
        x=$((x+1))
```

```
    done
```

```
    cp `printf "%s.mrc" "$dir"` "$dir"
```

```
    cd "$dir"
```

```
    echo "Processing images"
```

```
    find . -name "*.tif" -print0 | xargs -0 -l {} -n 1 -P 4 \
```

```
        mogrify -deskew 40 -gravity east \
```

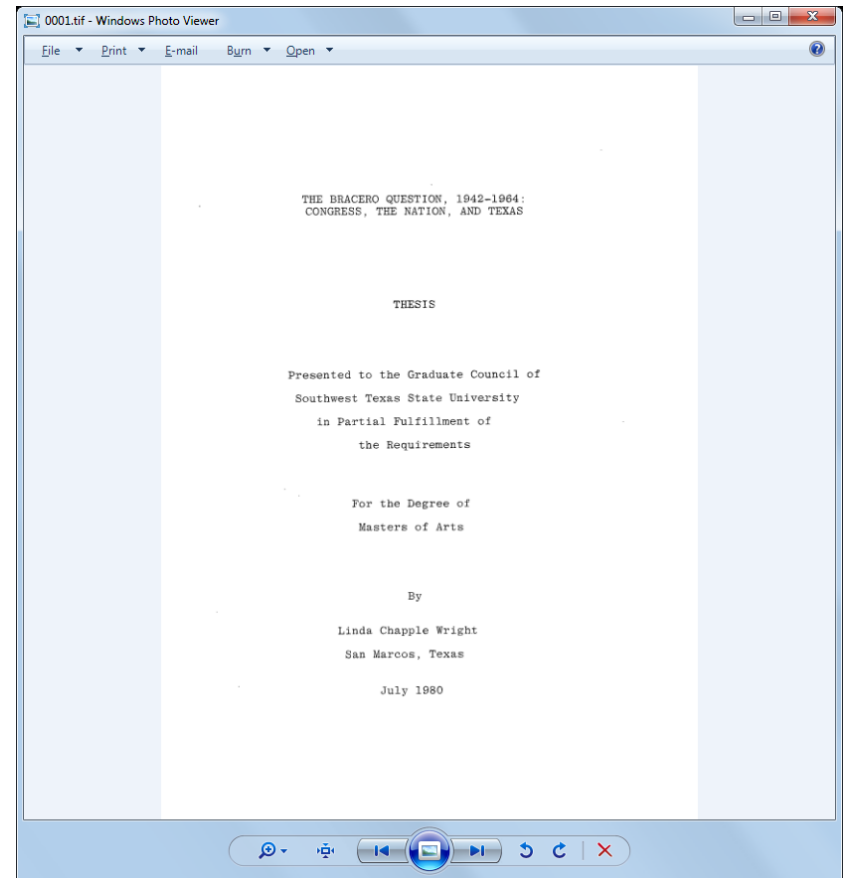
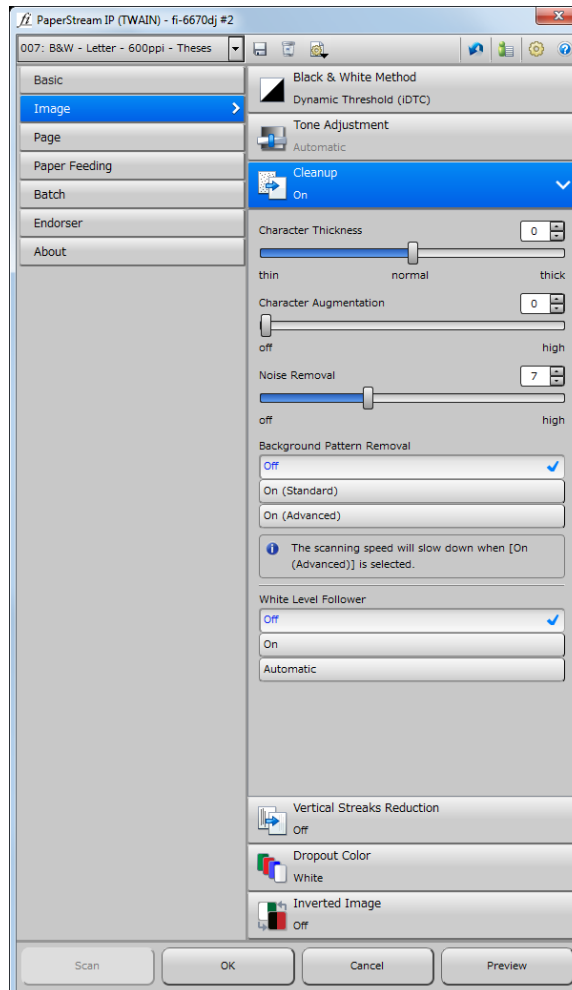
```
        -extent 5100x6600 -compress group4 "{}" 2>/dev/null
```

```
    echo "Images processed"
```

```
    cd ../..
```

```
done
```

Software Agnosticism

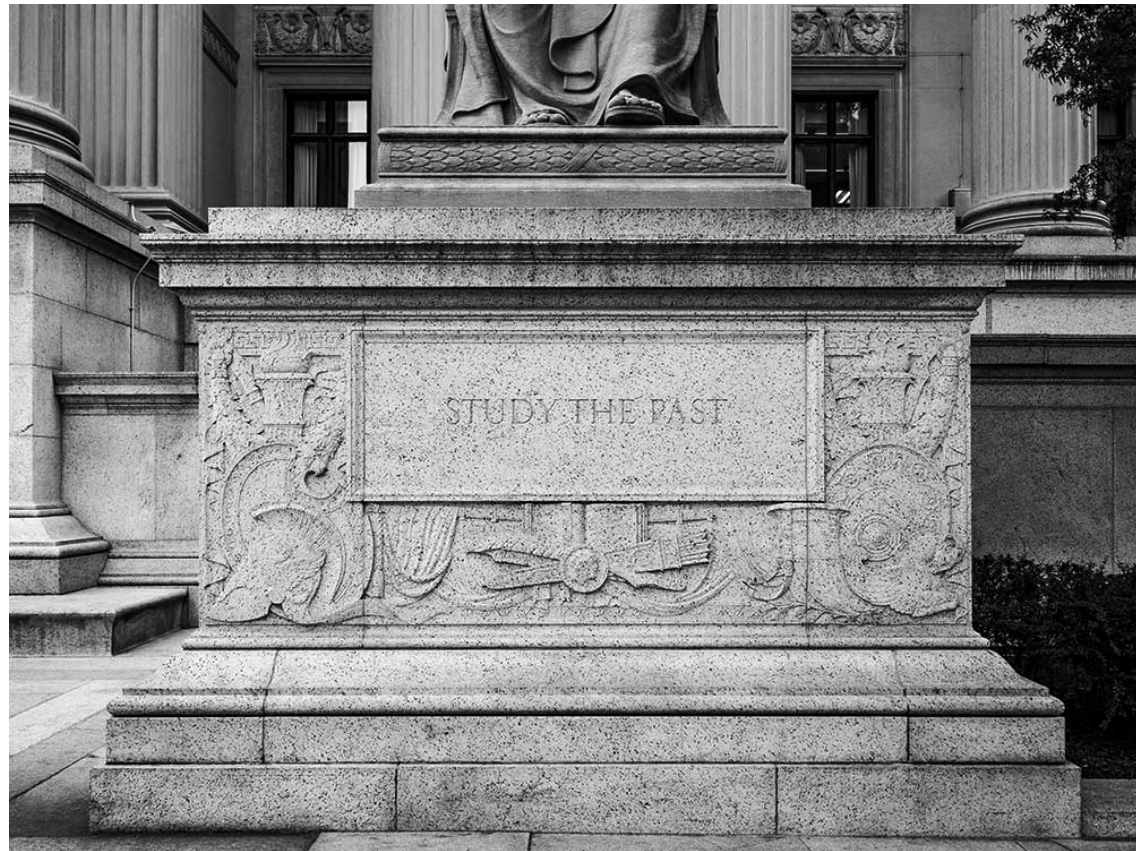


Cast of Characters

- ❖ Joan Heath, Assistant Vice President & University Librarian
- ❖ Kristine Toma, University Archivist
- ❖ Megan North, Assistant Archivist, University Archives
- ❖ Ray Uzwysyn, Director, Collections and Digital Services
- ❖ Jeanne Hazzard, Digital Collections Librarian
- ❖ Stephanie Towery, Copyright Officer
- ❖ Aaron Sinkar, Assistant Director, Administrative Services
- ❖ Misty Hopper, Head, Cataloging and Metadata Services
- ❖ Selene Hinojosa, Collection Development Librarian
- ❖ Paivi Rentz, Library System Coordinator
- ❖ Amy Eoff, Database Management Services Assistant
- ❖ Cliff Wood, Library Facilities Assistant
- ❖ Joshua Van Meter, Library Facilities Assistant
- ❖ Digital Imaging Technicians
 - Oscar Martinez, Project Lead
 - Luke Sebree
 - Ashton Woodward
 - Vincente Rangel
 - Wenlan Cai (now graduate research assistant with University Archives)

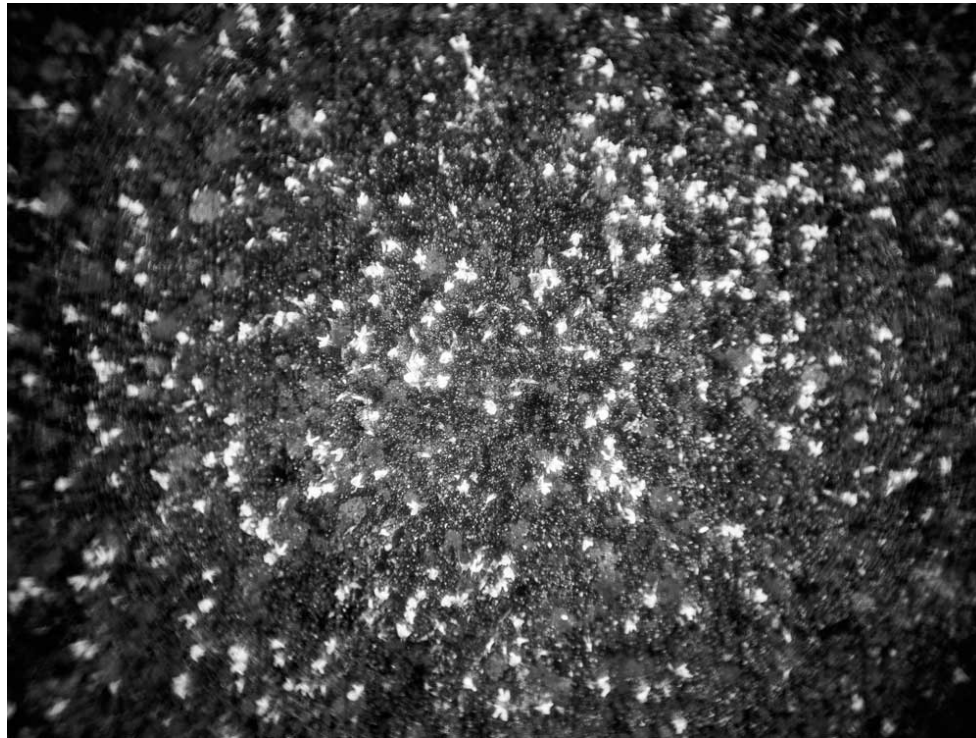
Project Management

- ❖ Fail quickly and often to find strong/weak links in the production chain



Training vs. Teaching

- ❖ Everything is the same until it isn't and everything is unique, except it's not.
 - **Binding, printing, paper quality changes with each thesis**



Fail fast = Skip problems

- ❖ No Ctrl+Z for most of the process



Leverage Strengths

❖ Human

- Kinesthetic
- Aesthetic
- Intuitive
- Creative
- Pattern recognition
- Reliable (with oversight)

❖ Computer

- Works weekends/holidays
- Repetitive
- Pattern recognition
- Reliable (with oversight)



Provide Training Wheels

- ❖ Start with looking and seeing: Quality Control
- ❖ PHP Web Application for sorting Color pages from B&W



Take them away!

- ❖ **Binding, printing, paper quality changes with each thesis**
- ❖ Personal paths to the end product, not a checklist of steps
- ❖ 'Problem' volumes are now opportunities for creativity and growth



- 0.toRescanBW
- 1.toPostProcess-BW
- 1.toPostProcess-Color
- 1a.Problems
- 2.toQC-BW
- 2.toQC-Color
- 3.QCing
- z.Batch2_toAccess
- z.Batch2_toPreservation

Quality Control for Documents



Questions?

