# EMOP

## EARLY MODERN OCR PROJECT

# eMOP guide to OCR

Flowcharting a Course Through Open-Source Waters

# eMOP Info

## eMOP Website

- emop.tamu.edu/

- TCDL 24x7 Presentation
  - emop.tamu.edu/news
- More information
  - emop.tamu.edu/TxDHC-flowcharts
- eMOP Workflows
  - emop.tamu.edu/workflows
- Mellon Grant Proposal
  - idhmc.tamu.edu/projects/Mellon/eMOPPublic.pdf

## More eMOP

- **Facebook**
  - Early Modern OCR Project

- **Twitter**
  - #emop
  - @IDHMC_Nexus
  - @matt_christy

# The Problems

**Early Modern Printing**

- Individual, hand-made typefaces

- Worn and broken type

- Poor quality equipment/paper

- Inconsistent line bases

- Unusual page layouts, decorative page elements,

- Special characters & ligatures

- Spelling variations

- Mixed typefaces and languages

The Maiden's Bloody Garland;
Or, HIGH-STREET TRAGEDY:

A PREFACE

...nes al our bleſſings, health, wealth, and proſperitie to the increaſe of Satans kingdome, are there abuſed: that not onſlie they are tearmed, as of late The ſchoole of abuſe, by one; The ſchoole of Bauderie by another; The neſt of the Diuel, and ſinke of al ſinne, by a third, ſo long agoe, The chaire of peſtilence, by Clement Alexandrinus; by Cyril, and Saluianus The pompe of the Diuel; the ſoueraigne place of Satan, by Tertullian.

And albe I cal them, A ſecond and third blaſt, &c. yet do I not ſo, as though there were no moe blaſtes, or dehortations fro them, or inuectiues againſt them beſide. For in al ages the moſt excellent men for learning haue condemned them by the force of eloquence, and power of Gods woorde (as I am to proue vpon anie good occaſion offered). But ſo do I tearme them

# The Problems

**Document/Image Quality**

- Torn and damaged pages
- Noise introduced to images of pages
- Skewed pages
- Warped pages
- Missing pages
- Inverted pages
- Incorrect metadata
- Extremely low quality TIFFs (~50K)

TCDL : eMOP: Flowcharting a Course Through Open-Source Waters

gretie abaſt, ot enſuarble ſappoſing ſurelp þat
thoſe habeſte æ bampoe . anoe be cauſe alſo thoſe
habeſte no þype ne remedp be fore the tyme of the
tole . ... mentes of the cherch . Soothle be cauſe
e ſenſe ... other theſe than ... ... p am glade
... ... ... þoolee þere come ... te on what theſe
... fo anoe ſcappopſte eternal oamnapeoo
... . Slable what ſum euer thoole before
... the telle

... ... goldſmpth tolee the noble gyl pur
... ... ... hpeoz ſooenlp ... ... ... ſaupoe
... ...

p ... ... ... ... ... ... oſpoe me on my leſ
... ... ... ... ... on the þorlee ſe touch þa
gþe ... ... e oppn ... ... jepbre . Alſo p conſ
tynooplt on the foole ... ue of wronkpnnes . on to
me laſt ... ... of an engl cuſtome . Aenerthpleo het
we not ... ... ... ... ... hpt oeſpleſou me
... ... ... ... ... ſar e oooee no ... que that
... ... ... then tymes p roſe ageynſt mp ſelfe
... ... ng ... ne τ caſte aoeep the foole
... ... ... ... ... ... ... ... Butano
... ... ... ... ... ... ... ... ...
... ther ... ... θ. pleno pthapnte to begn
ke aſ ... ... oneo of mpne olt cuſtome . oehehh's
waſ ouer came . τ oτ ... ... awgne boloe pat to take
τ cuſtome of the ſame ſinne that was ... ſame of the
onneſaoulle taking τ appetite Λ noue amonge this

# Wrangling Data

## The Numbers

- **EEBO**: ~125,000 documents, ~13 million pages images (1475-1700)

- **ECCO**: ~182,000 documents, ~32 million page images (1700-1800)

- **TCP**: ~46,000 double-keyed hand transcriptions (44,000 EEBO, 2,200 ECCO) - Groundtruth

- **Total**: >300,000 documents & 45 million page images.

## The Data

- **ECCO page images** (1 pg/image)

- **ECCO original OCR results** (doc-level XML files)

- **ECCO TCP transcriptions** (doc-level XML and text files)

- **EEBO page images** (2 pgs/image)

- **EEBO TCP transcriptions** (doc-level XML and text files)

Wrangling Data



ECCO
Eighteenth Century Collections Online

Page images and GaleXML OCR transcripts ~ **182,000 docs**

TCP
Text Creation Partnership

hand-keyed transcriptions: **44,000** EEBO, **2200** ECCO

Groundtruth

XML

XML

Convert doc-level XML to page-level XML and Text

Mine XML files for metadata

XML Text

45 million Total Page Images

IDHMC NAS

file locations

eMOP DB

Page images and metadata ~ **125,000 docs**

Mine metadata for DB ingestion

EEBO
Early English Books Online

eMOP Tesseract Training & OCR'ing Workflow

# Training Tesseract

**Aletheia**

Created by PRImA Research Labs. A team of undergraduates uses Aletheia to identify each glyph on the page images, and ensure that the correct Unicode value is assigned to each. Aletheia outputs an XML file containing all identified glyphs on a page with their corresponding coordinates and Unicode values.

# Training Tesseract

**Franken+**

1. Takes Aletheia's output files as input.
2. Groups all glyphs with the same Unicode values into one window for comparison.
3. Mistakenly coded glyphs are easily identified and re-coded.
4. A user can quickly compare all exemplars of a glyph and choose just the best subset, if desired.
5. Uses all selected glyphs to create a Franken-page image (TIFF) using a selected text as a base.
6. Outputs the same box files and TIFF images that Tesseract's first stage of native training.
7. Also allows users to complete Tesseract training using newly created box/TIFF file pairs, and add optional dictionary and other files.
8. Outputs a .traineddata file used by Tesseract when OCRing page images.

# Control Tools

- **Query Builder**: Users can build queries to find specific documents, or sets of documents, in the eMOP DB. Document sets can be labeled or grouped with an identifier.

- **Data Downloader**: Available via the QB. Allows identified items to be downloaded: page images, OCR results, associated groundtruth.

- **eMOP Dashboard**: Users select documents from the database, individually or with the identifier created in the QB, for OCRing with specific typeface training. The Dashboard displays results with Juxta/RETAS scores when groundtruth is available.



**OCR Process**

eMOP Dashboard

**Data Control**

Query Builder / Data Downloader

.tif, .xml, ..txt

eMOP Database

IDHMC NAS

emop_ controller

# emop_controller

A java program that runs on the Brazos High Performance Computing Cluster and ensures maximum utilization of the 128 processors available for our use by scheduling various functions and processes of the controller separately.

1. The selected documents are marked in DB as scheduled and the documents are queued for processing.

2. A cron job continuously checks the scheduled queue and OCRs any unscheduled pages.

3. Each page is OCRd with Tesseract.

4. After OCRing each page's hOCR output file is de-noised.

5. hOCR pages that have matching Groundtruth are scored using Juxta and RETAS algorithms.

6. File paths and scores are written to the eMOP DB.

7. hOCR documents are examined again and given an estimated correctability (ECORR) score.

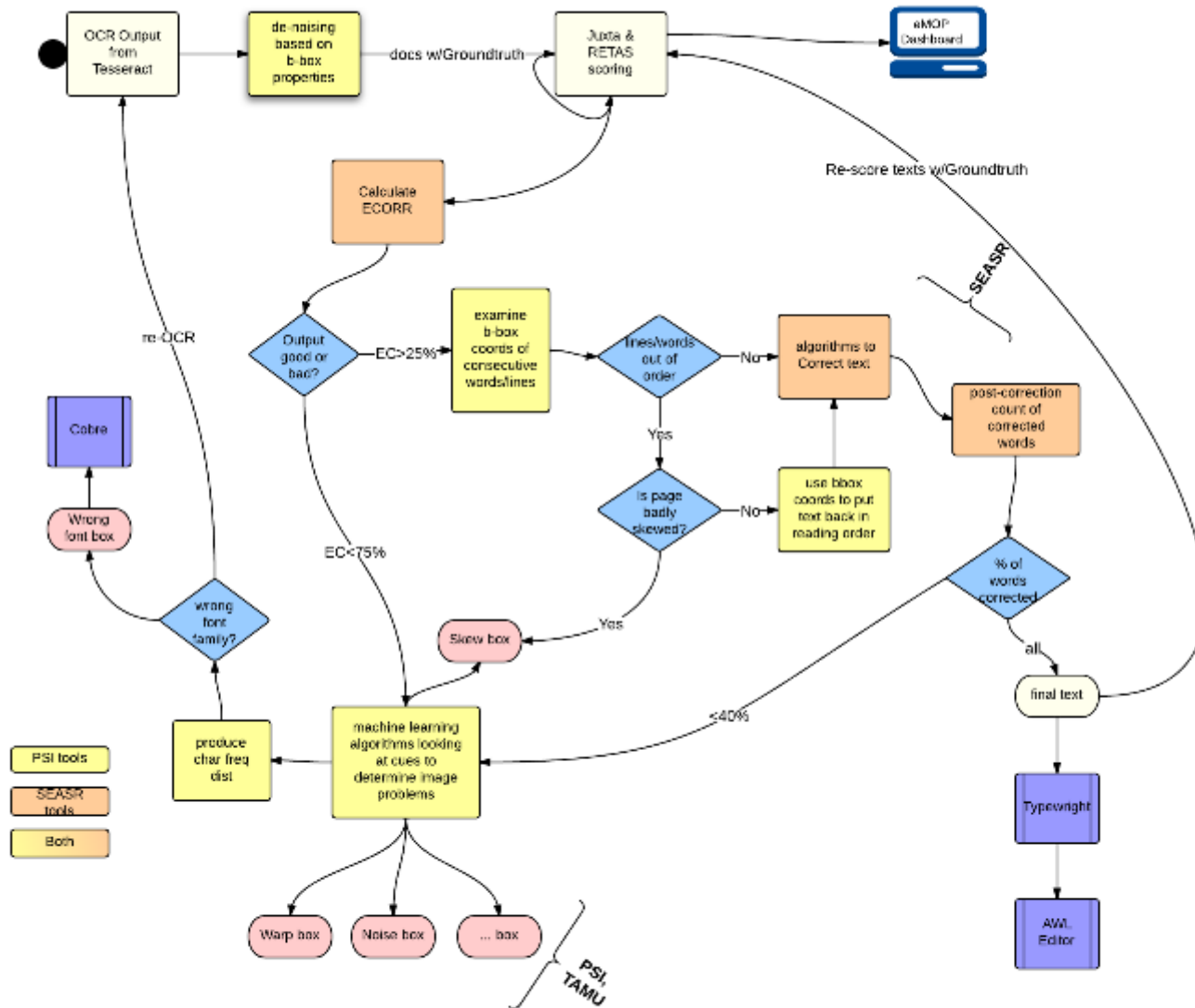8. hOCR page results are sent to post-processing triage.

emop_controller

# De-noising

De-noising

Post-Processing Triage

# Post-Processing Triage

1. De-noise hOCR results
2. Juxta/RETAS scores for docs w/Groundtruth
3. Calculate ECORR

ECORR > 25%

ECORR < 75%

### Attempt to further correct text (SEASR)

4. Check bbox coords of *consecutive* lines/words to identify if any are out of order.
    4.a. If so, check to see if the page is too badly skewed.
        4.a.i. If so, mark page as skewed in the eMOP DB
        4.a.ii. Otherwise, use bbox coords to make consecutive boxes adjacent.
    4.b. Otherwise, send page for correction.

5. Correct page's text as much as possible.

6. Count corrected words.

7. Compute ratio of corrected words to total words on page.
    7.a. If ratio is less than 40%, send to PSI queue for analysis.
    7.b. Send all processed pages to TypeWright for crowd-sourced hand correction.

### Analyze hOCR elements to diagnose page image problems (PSI)

4. Machine-learning algorithms examine OCR results to determine what is wrong with the page image: skewed, warped, noisy, wrong font, etc.

5. Create a character frequency distribution for pages diagnosed with wrong font.

6. Does the character frequency distribution indicate the wrong typeface family was used: i.e. blackletter/roman/italic?
    6.a. If so, send back for re-OCRing with different font family.
    6.b. Otherwise, or if this page has already been OCRd with multiple typeface families, send to Cobre for typeface identification by an expert.

# Post-Processing Triage

correct-ability score

>25% → SEASR corrections

<75%

PSI analysis and categorization of document problem

<40%

% of words corrected

all

DB of page images with "problems" tags

TypeWright

# TypeWright

- [www.18thconnect.org/typewright/](http://www.18thconnect.org/typewright/)

- Available through [18thconnect.org](http://18thconnect.org)
  An aggregator of 18th Century metadata & content

**This will be the first time that EEBO documents are available as searchable text, and that they can be corrected by scholars for their own use.**

- Crowd-sourced correction tool for OCR.

- Currently available for ECCO collection, based on original ECCO OCR results.

- Will be ingested with eMOP OCR for ECCO and EEBO by year-end.

- Users can correct a whole document and then receive the text or a lightly encoded TEI XML version for use in a digital edition or whatever.

TypeWright



Edit button identifies "TypeWright" enabled Texts.

# Tags DB

- All page images that don't produce OCR text that meets our standards will be tagged in the eMOP DB with tags that indicate what is wrong with the page image.
  - skew, warp, noise, wrong font, etc.

- Tags can then be used to later apply the appropriate pre-processing techniques and then re-OCR the page to produce better results.

- **This will be the first time that any sort of comprehensive analysis has been done on the page images of these collections.**

- We'll end up (ideally) with a list of page images that simply need to be re-scanned.

# The end

For eMOP questions please contact me at :

mchristy@tamu.edu