



**IDHMC**  
TEXAS A&M UNIVERSITY

**Initiative for Digital Humanities, Media, and Culture**

# **Early Modern OCR Project (eMOP) at Texas A&M University: Using Aletheia to Train Tesseract**

Katayoun Torabi, PhD Student in English

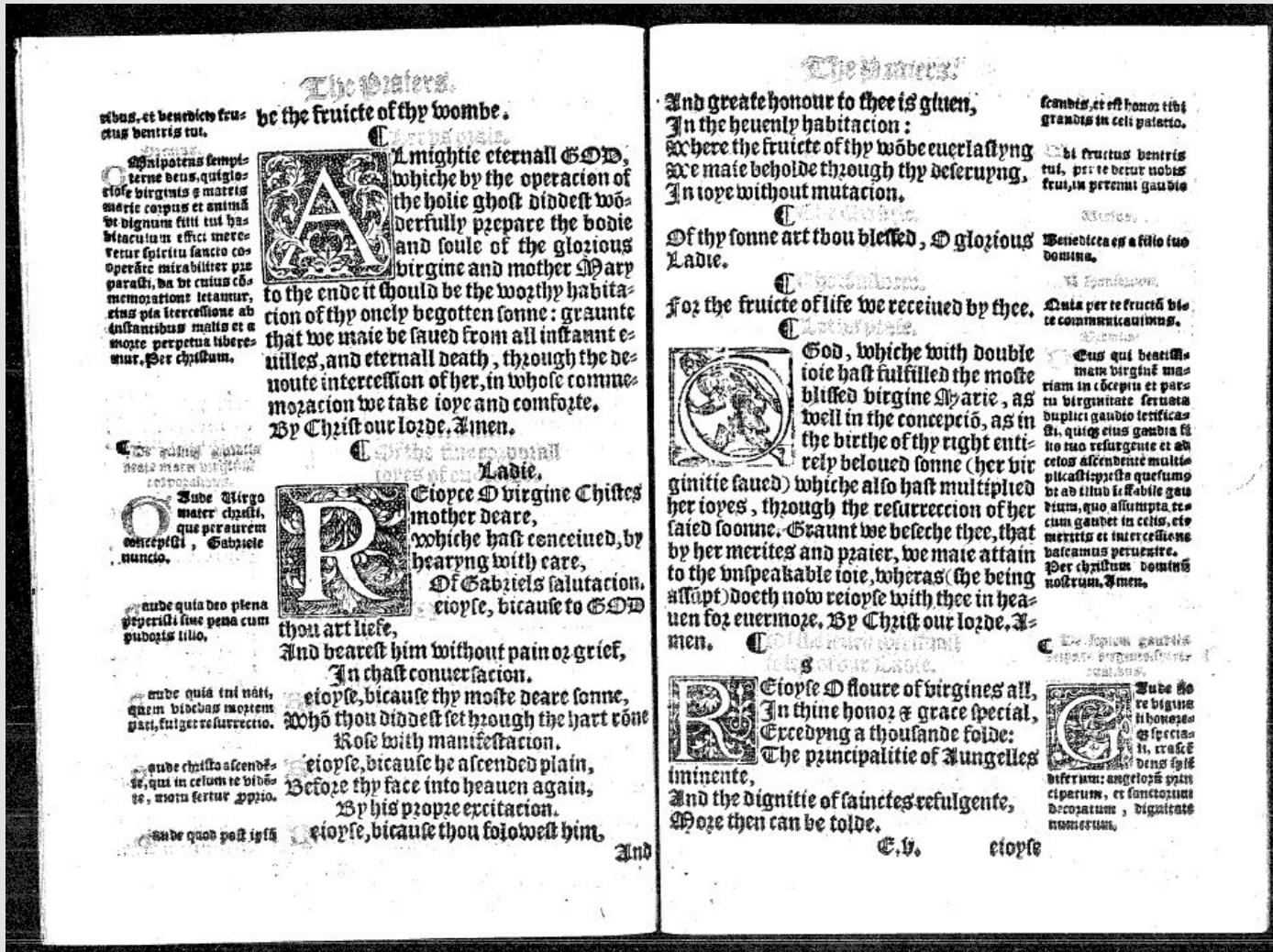
Bryan Tarpley, PhD Student in English

Dr. Jessica Durgen, Post-Doctoral Fellow (not present)

# The Problem:

- The IDHMC is tasked with performing OCR on over 45,000,000 page images of historical documents which were scanned using poor resolution scanners with problems such as idiosyncratic fonts and bleed-through.
- In-house, we have at our disposal 2 programmers and ~ 5 part time students in the English department and the Digital Humanities. We are also privileged to work with various collaborators, including PRImA.
- We may only work with open-source solutions (no ABBY FineReader).
- The time frame for our project is approximately 2 years.

# Historical Documents



# Ligatures



# Italics



*Germany.*

# Blackletter

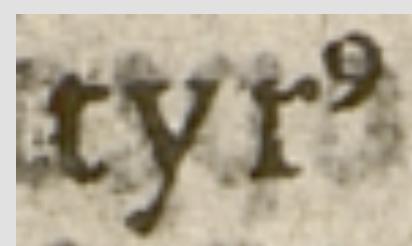
96

THE TREASVRE

And because it hath much heat, it expelleth what soever is vncleane, and therfore restoreth the natural heat . For age is nothinge els then a lessening and diminishing of natural heat , which is therfore diminished, because the mouinge is hindred: as I saide of fire: for there is like reason of this vnto that. Mouing is hindred , because of þ abundance of earthly matter, because the earthe only hath very muche matter, and is destitute of all mouing. Therfore that water being of so temperate a heate, it shall nether vexe the hearte, nor noysom to the liuer: wherfore this onlye can perform the things that we haue spoken. When burning water therfore or Aqua vita reteininge his purity and subtil matter : by the longe mocyon of circulacio, hath put of and rid away his heat and Sharpnes of taste and smell , it is becomed of the

# Other Special Characters

- Other special characters encountered are letters with suspension marks, rotunda r's, and letters with superscripts



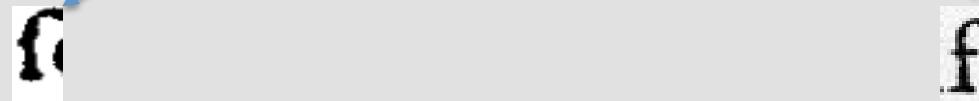
# Long ‘s’ in the word ‘greatnesse’



A close-up view of the word "greatnesse" in a black, serif font. The letter "s" is a long s, which is a character where the vertical stem of the 's' extends upwards from the top loop. The letters are set against a white background with faint gray grid lines, and a small portion of a green object is visible at the bottom left.

# Long 's' and 'f'

So high or low, dost raise thy formall hat:



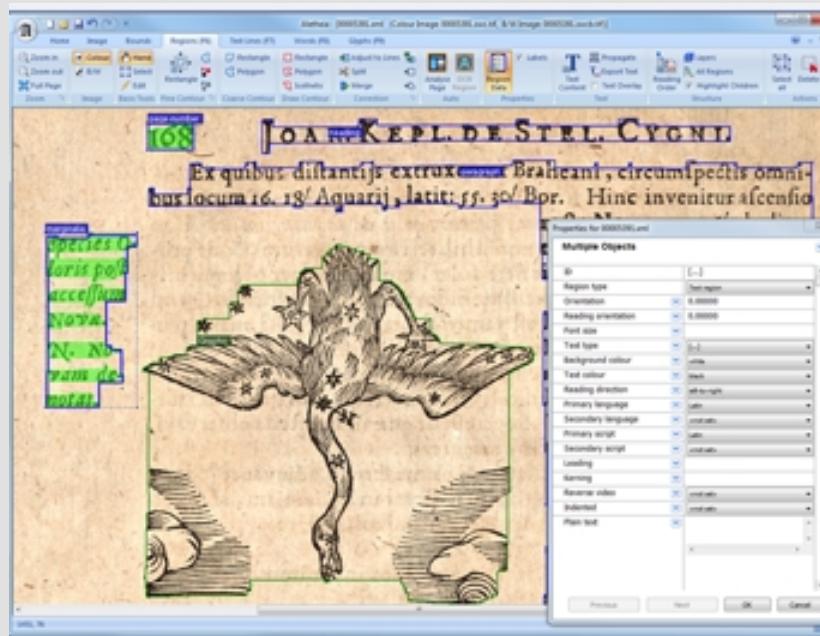
f

f

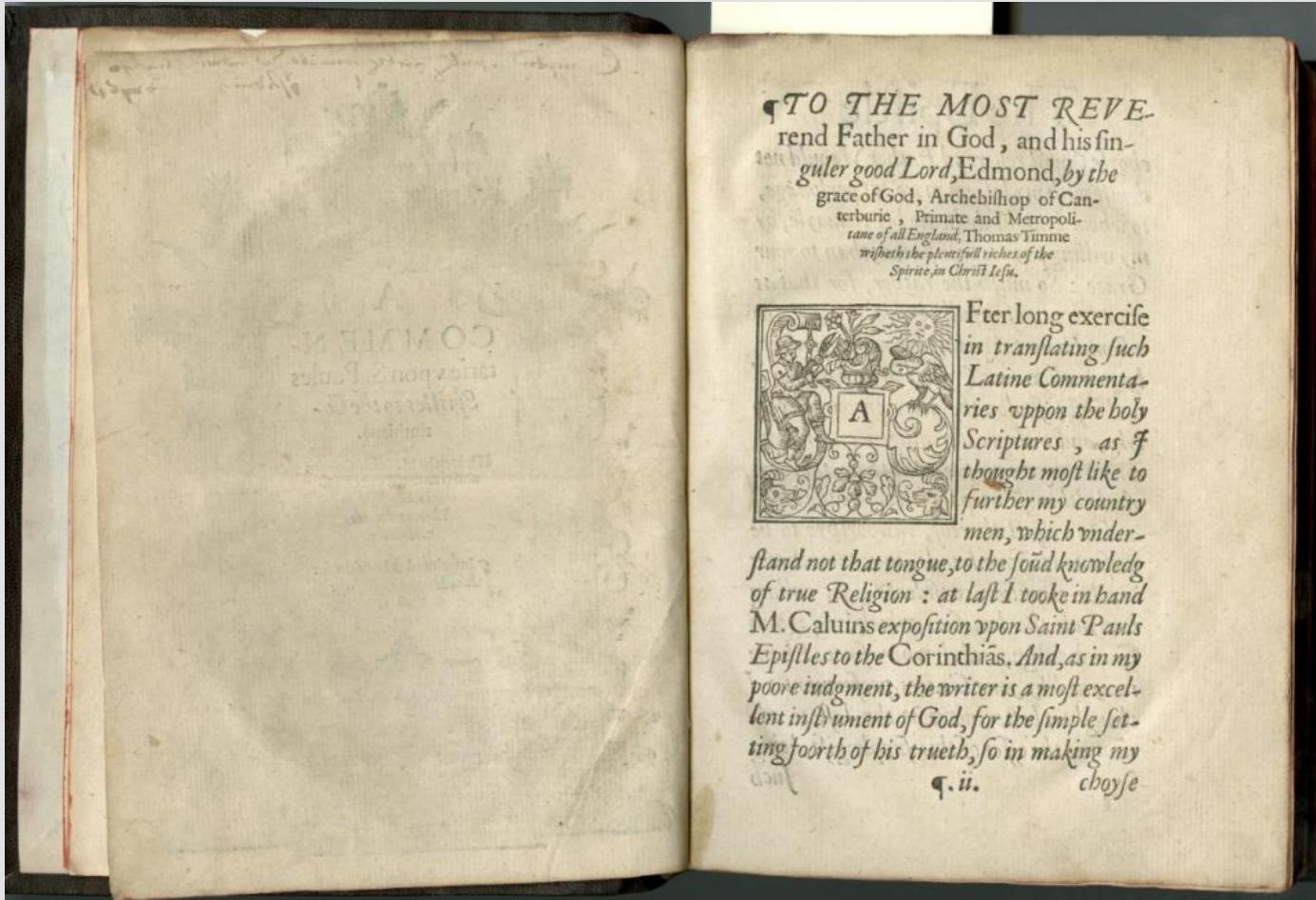
's' with a half cross-bar

'f' with a full cross-bar

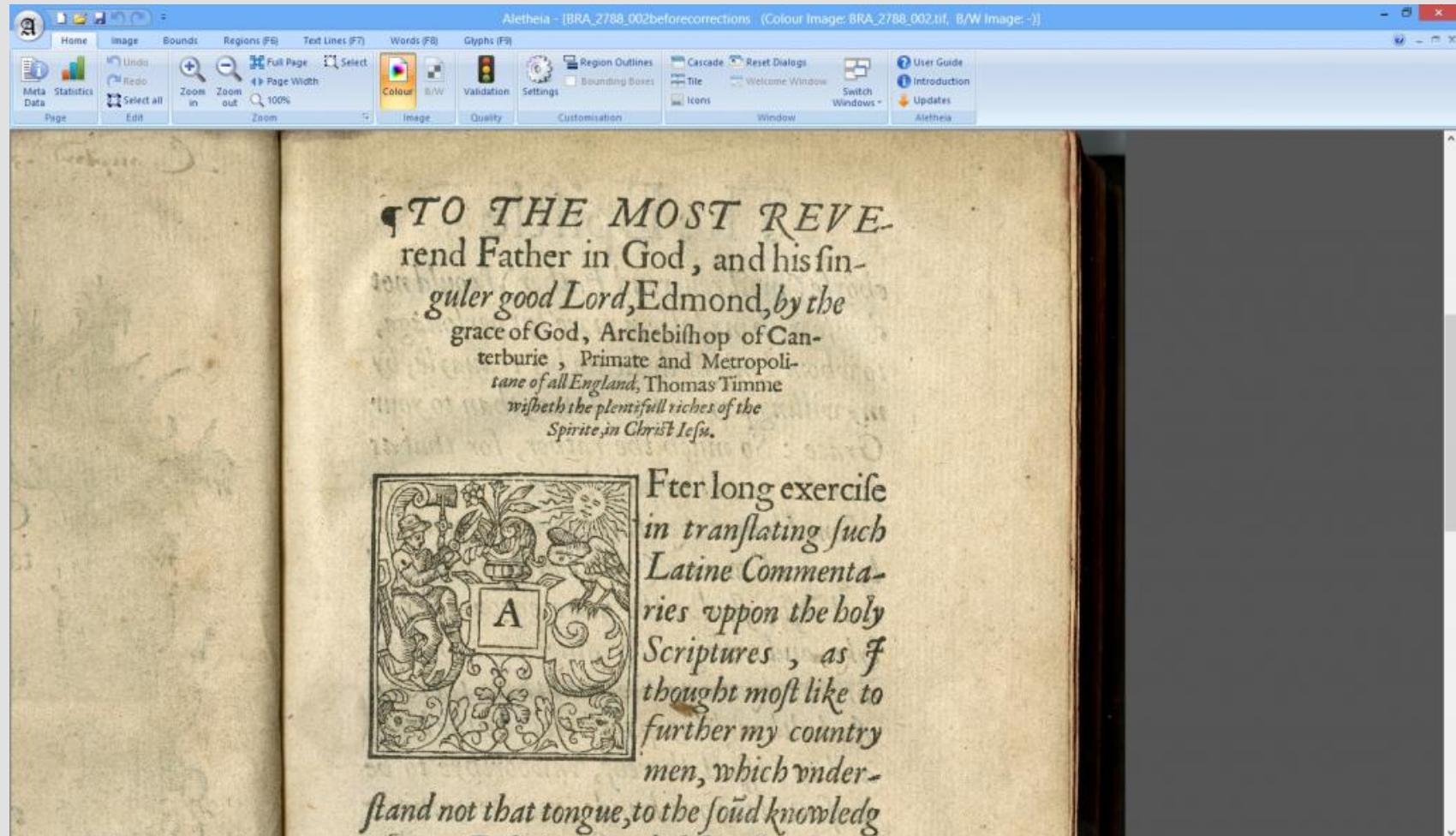
## PRImA, the developers of Aletheia Desktop



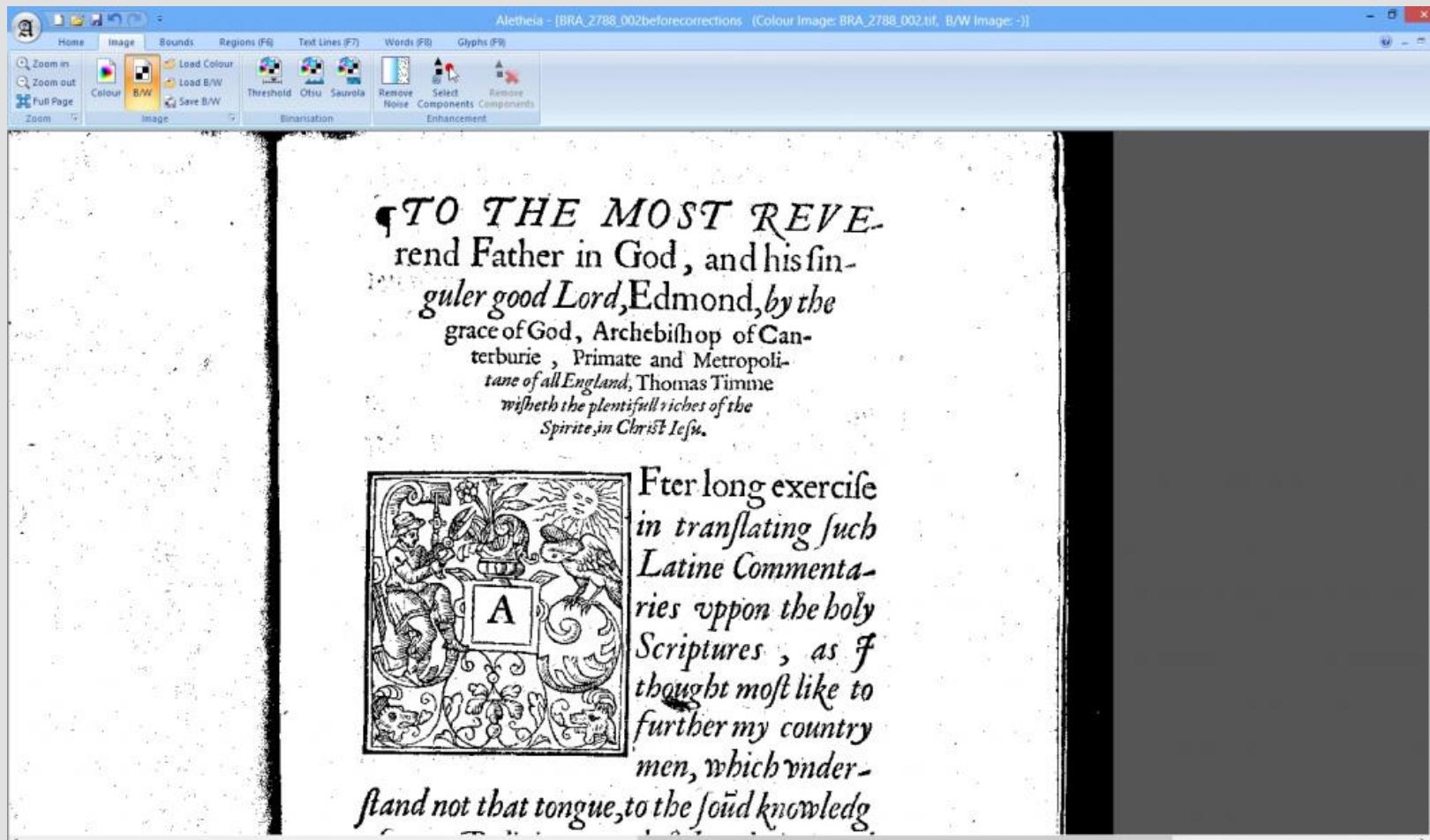
# TIFF Page Image



## Step 1: Upload the TIFF image in Aletheia

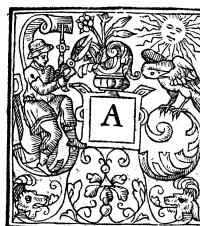


# Step 2: Binarize the image



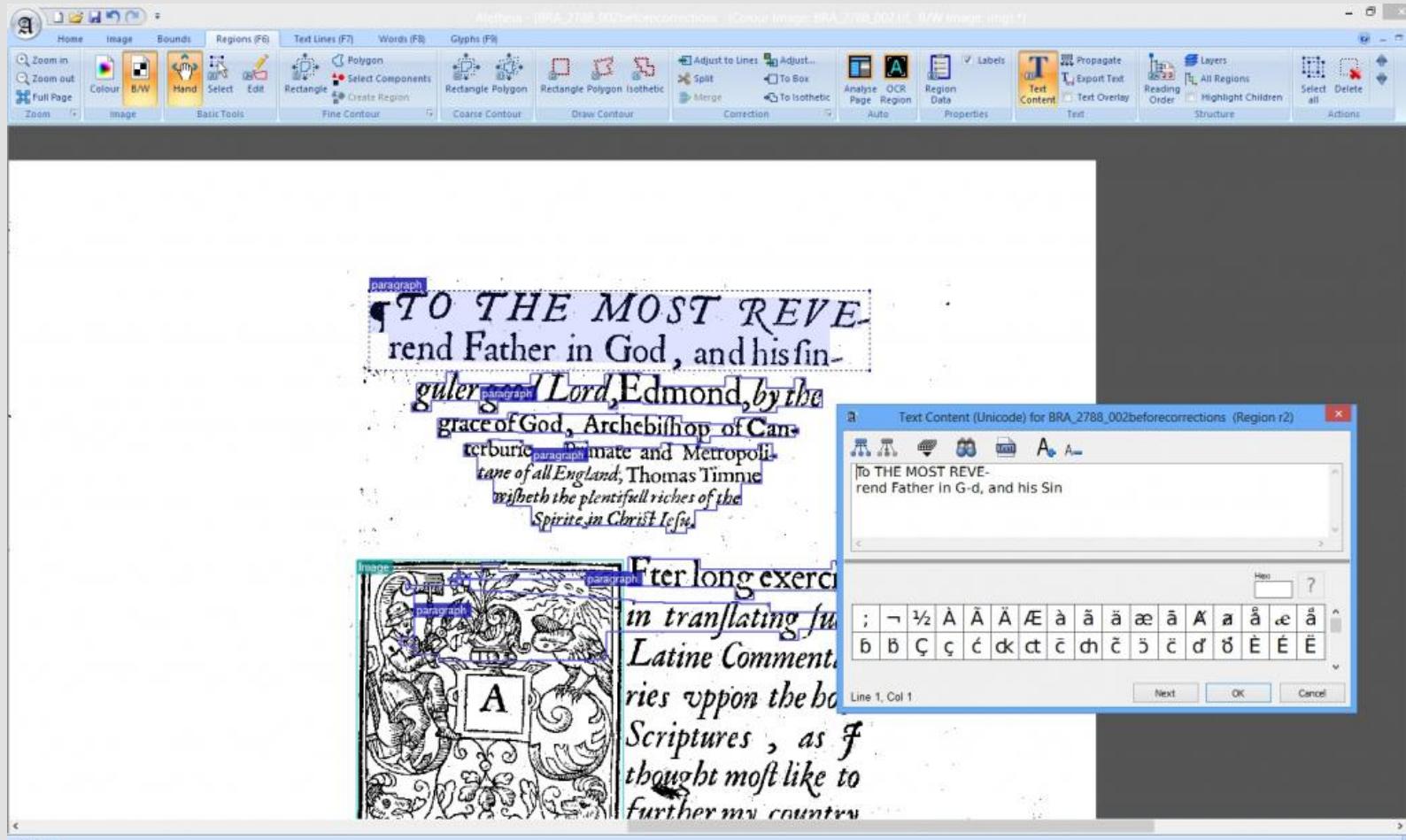
# Step 3: Remove noise

**T**O THE MOST REVER-  
 rend Father in God, and his singular good Lord, Edmond, by the grace of God, Archebishop of Canterbury, Primate and Metropolitane of all England, Thomas Tymme wibeth the plentifull riches of the Spirit in Christ Iesu.

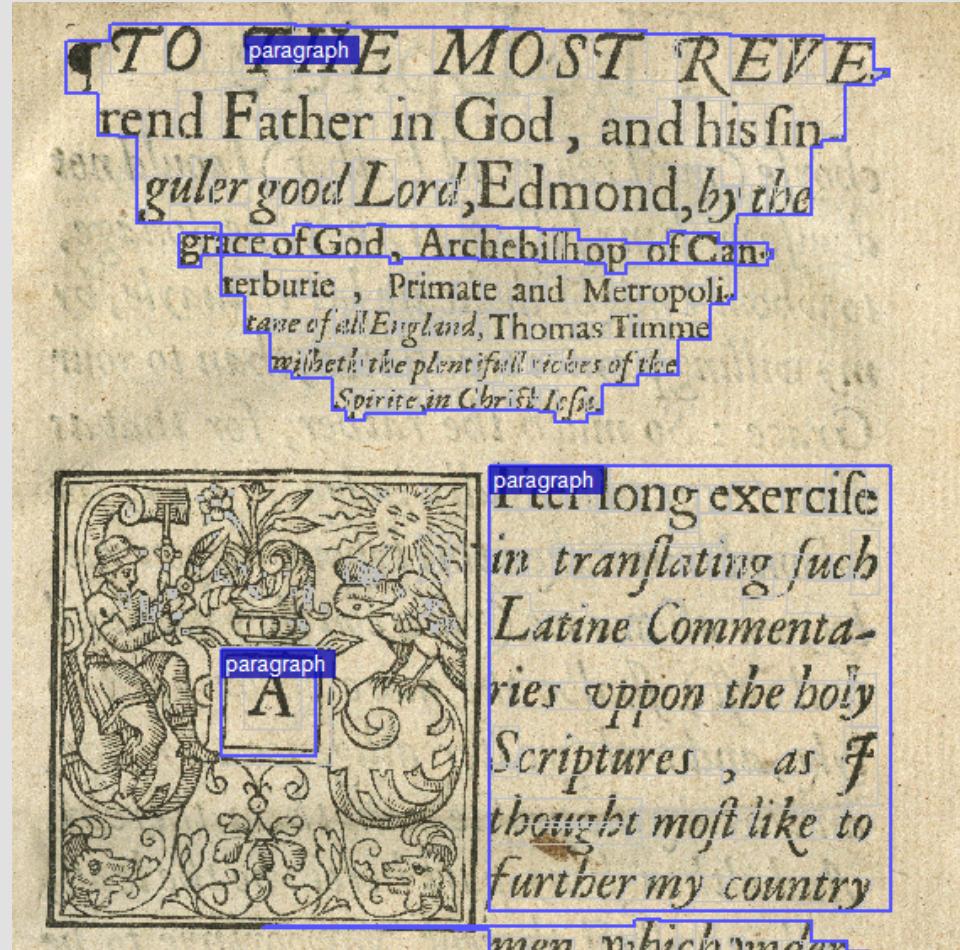


After long exercise in translating such Latine Commentaries vpon the holy Scriptures, as I thought most like to further my country men, which understand not that tongue, to the sond knowledg of true Religion: at laſt I tooke in hand M. Caluins expositiōn vpon Saint Pauls Epistles to the Corinthiās. And, as in my poore iudgment, the writer is a moſt excellent iſtument of God, for the ſimple ſetting forth of his truthe, ſo in making my ¶. ii. choyſe

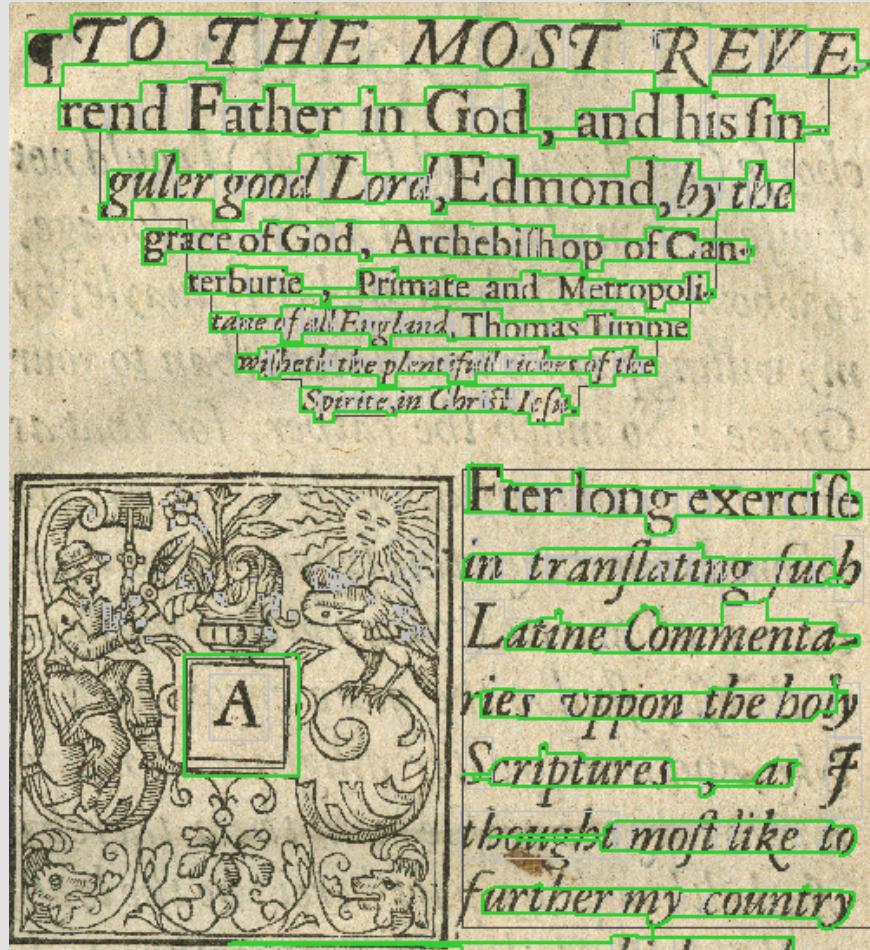
## Step 4: Running the automatic segmentation tool in Aletheia



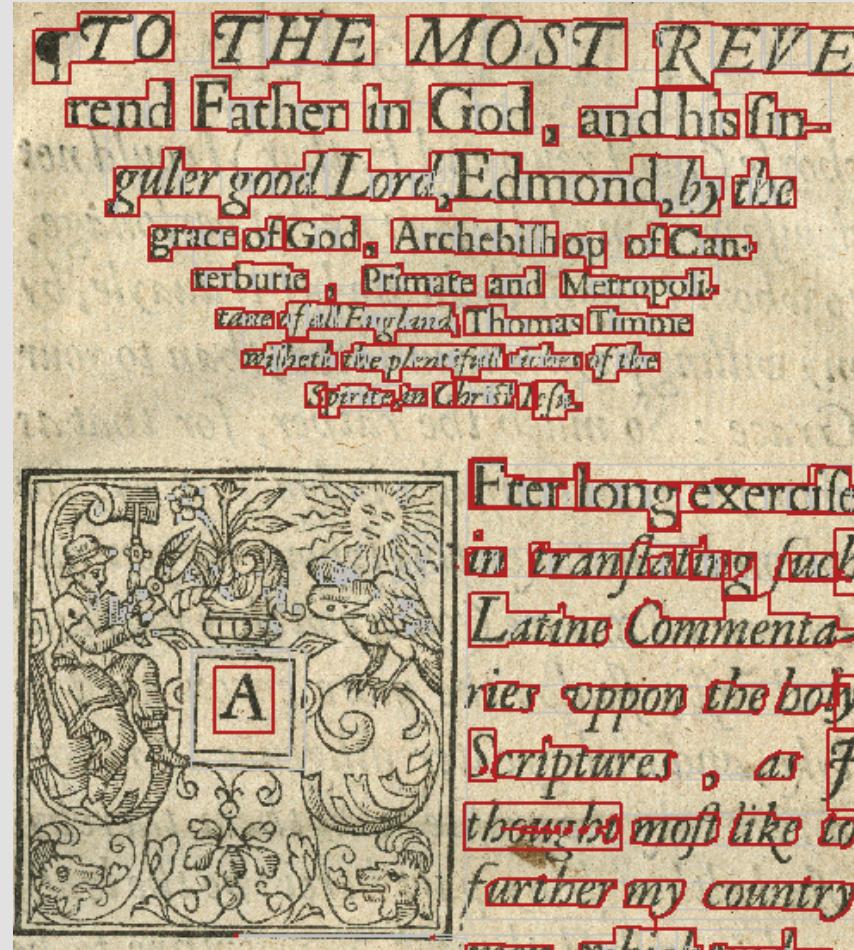
# Regions



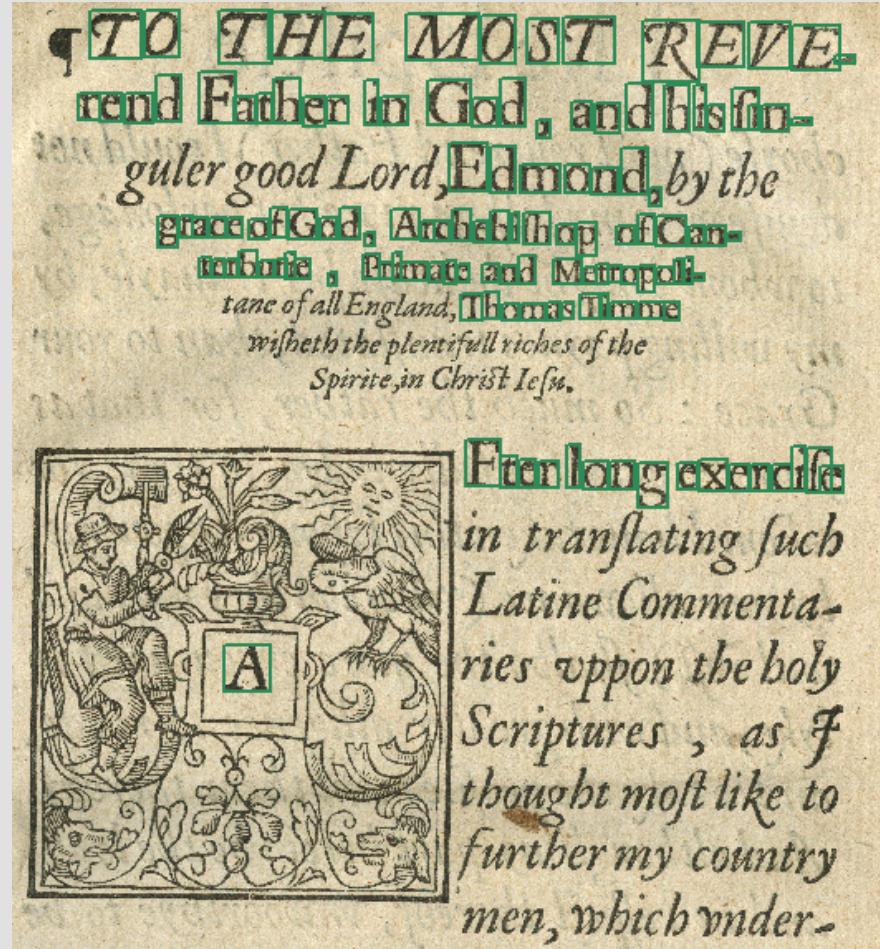
# Lines



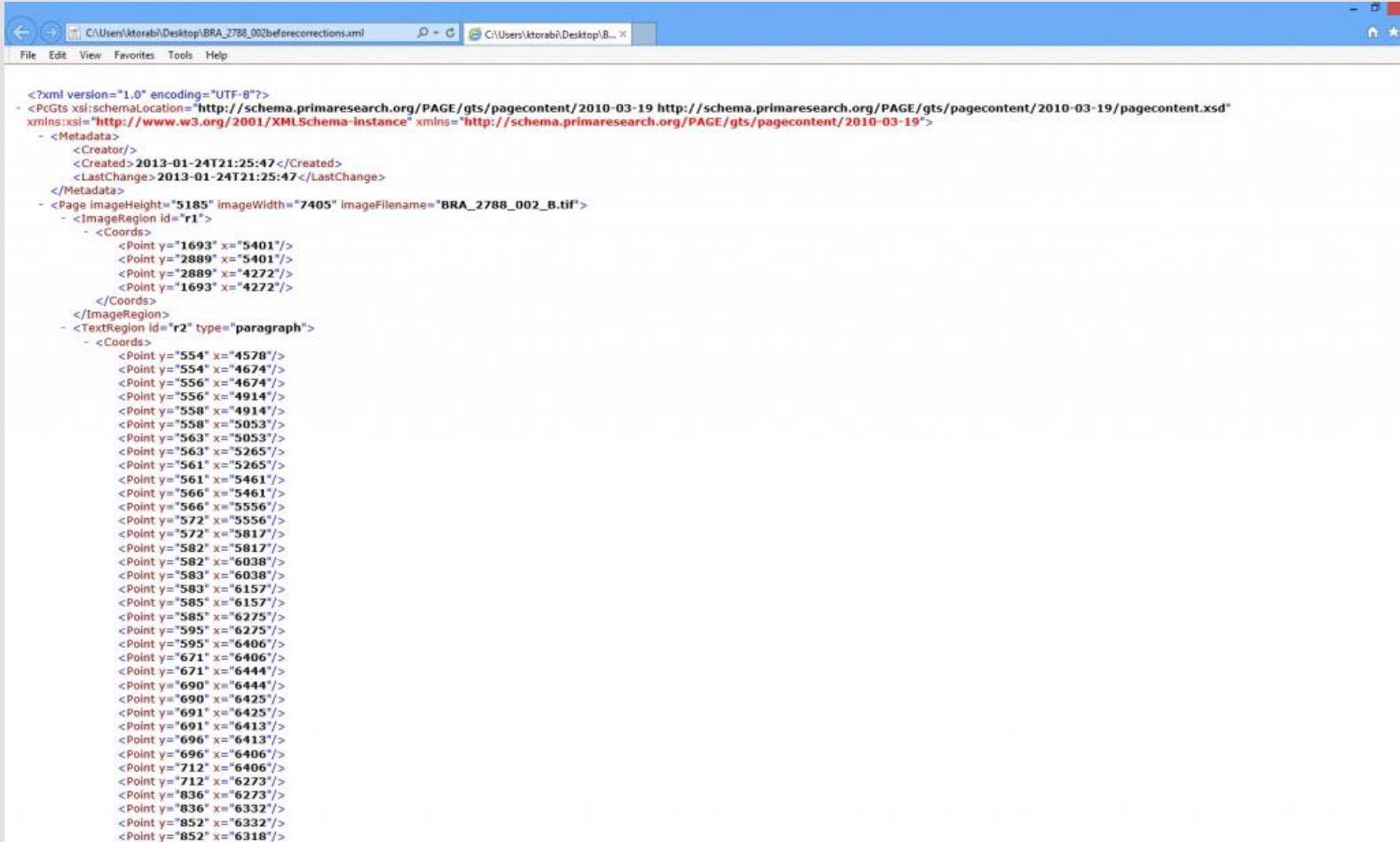
# Words



# Glyphs



The segmentation tool generates text for each character and a corresponding XML file

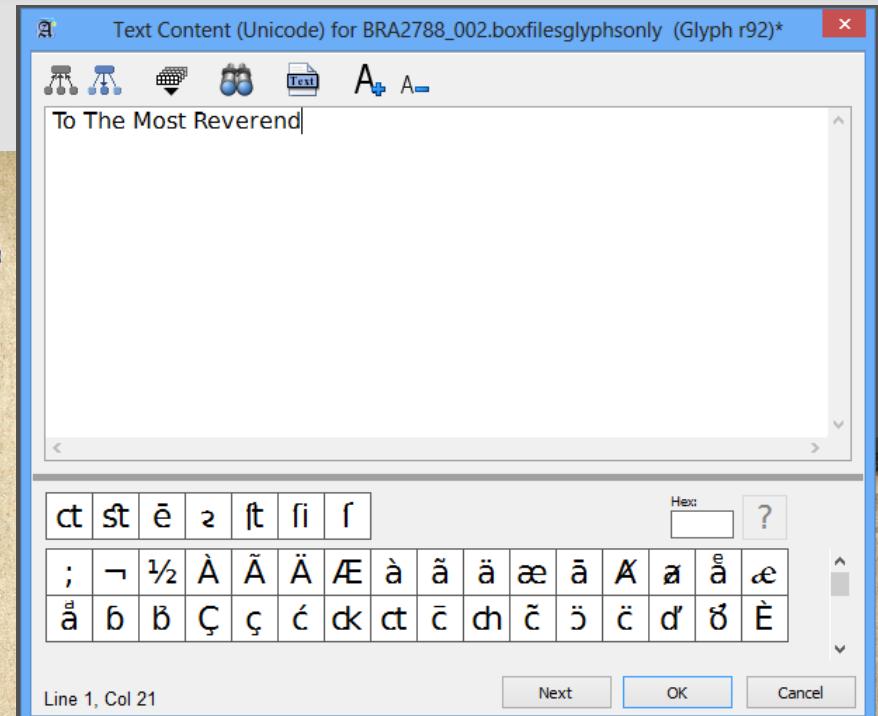
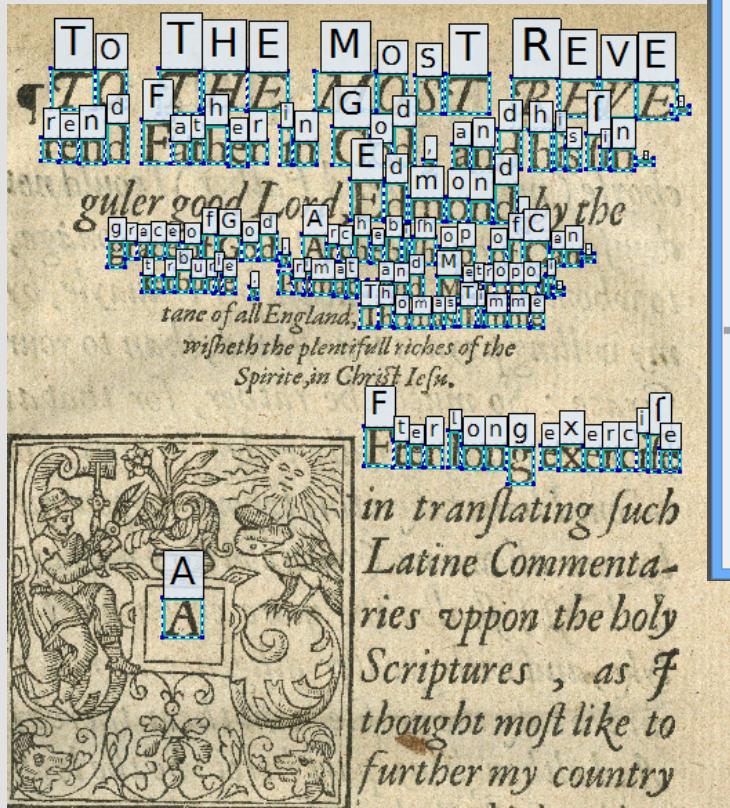


```

<?xml version="1.0" encoding="UTF-8"?>
<PG>
  xsi:schemaLocation="http://schema.primaresearch.org/PAGE/gts/pagecontent/2010-03-19 http://schema.primaresearch.org/PAGE/gts/pagecontent/2010-03-19/pagecontent.xsd"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xmlns="http://schema.primaresearch.org/PAGE/gts/pagecontent/2010-03-19">
    <Metadata>
      <Creator/>
      <Created>2013-01-24T21:25:47</Created>
      <LastChange>2013-01-24T21:25:47</LastChange>
    </Metadata>
    <Page imageHeight="5185" imageWidth="7405" imageFilename="BRA_2788_002_B.tif">
      <ImageRegion id="r1">
        <Coords>
          <Point y="1693" x="5401"/>
          <Point y="2889" x="5401"/>
          <Point y="2889" x="4272"/>
          <Point y="1693" x="4272"/>
        </Coords>
      </ImageRegion>
      <TextRegion id="r2" type="paragraph">
        <Coords>
          <Point y="554" x="4578"/>
          <Point y="554" x="4674"/>
          <Point y="556" x="4674"/>
          <Point y="556" x="4914"/>
          <Point y="558" x="4914"/>
          <Point y="558" x="5053"/>
          <Point y="563" x="5053"/>
          <Point y="563" x="5265"/>
          <Point y="561" x="5265"/>
          <Point y="561" x="5461"/>
          <Point y="566" x="5461"/>
          <Point y="566" x="5556"/>
          <Point y="572" x="5556"/>
          <Point y="572" x="5817"/>
          <Point y="582" x="5817"/>
          <Point y="582" x="6038"/>
          <Point y="583" x="6038"/>
          <Point y="583" x="6157"/>
          <Point y="585" x="6157"/>
          <Point y="585" x="6275"/>
          <Point y="595" x="6275"/>
          <Point y="595" x="6406"/>
          <Point y="671" x="6406"/>
          <Point y="671" x="6444"/>
          <Point y="690" x="6444"/>
          <Point y="690" x="6425"/>
          <Point y="691" x="6425"/>
          <Point y="691" x="6413"/>
          <Point y="696" x="6413"/>
          <Point y="696" x="6406"/>
          <Point y="712" x="6406"/>
          <Point y="712" x="6273"/>
          <Point y="836" x="6273"/>
          <Point y="836" x="6332"/>
          <Point y="852" x="6332"/>
          <Point y="852" x="6318"/>
        </Coords>
      </TextRegion>
    </Page>
  </PG>

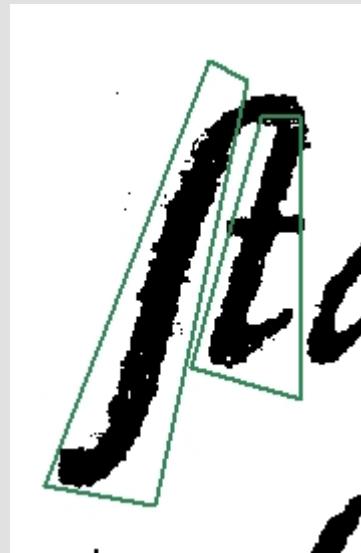
```

# Text



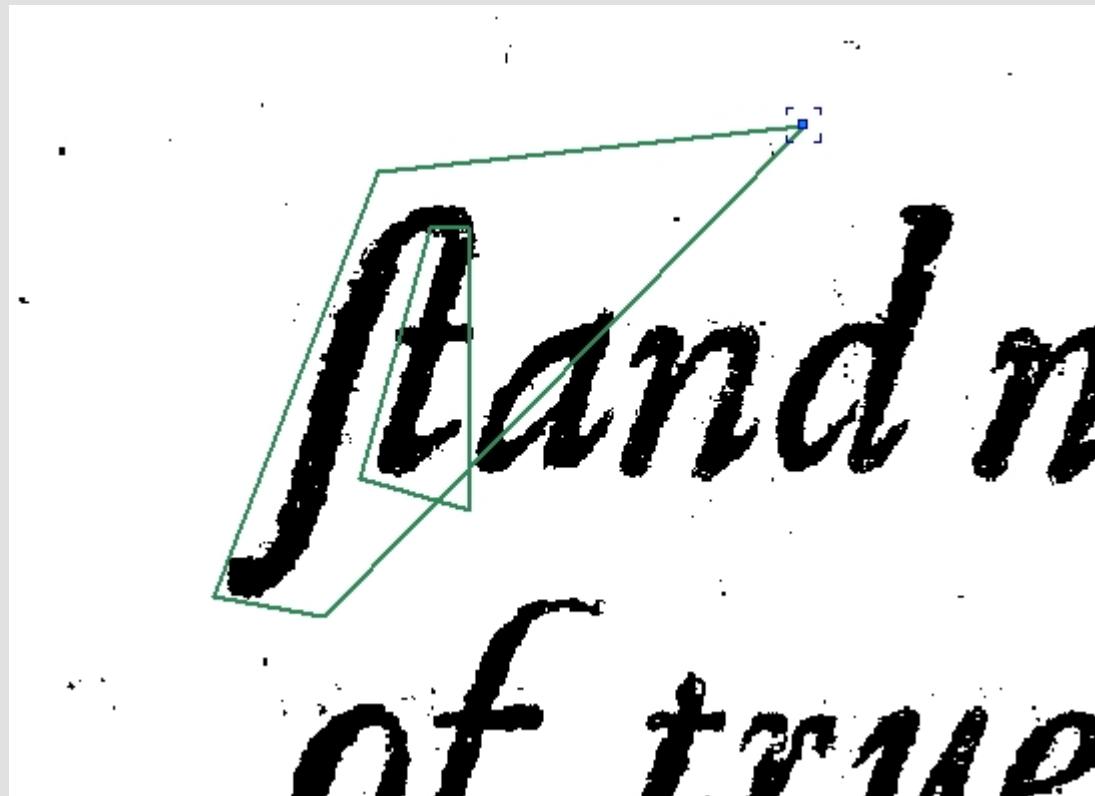
# Misreads

- Ligature boundaries are sometimes confused



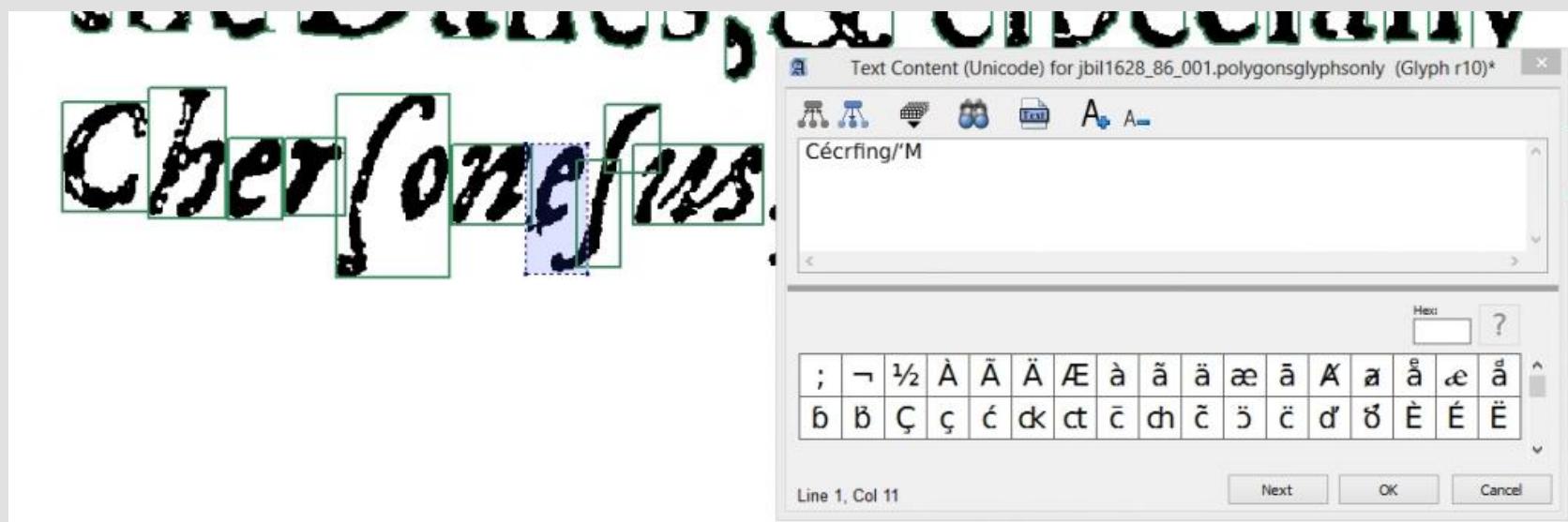
## Step 5: Glyph Correction

- The parameters of a glyph may be corrected using the edit button to adjust lines



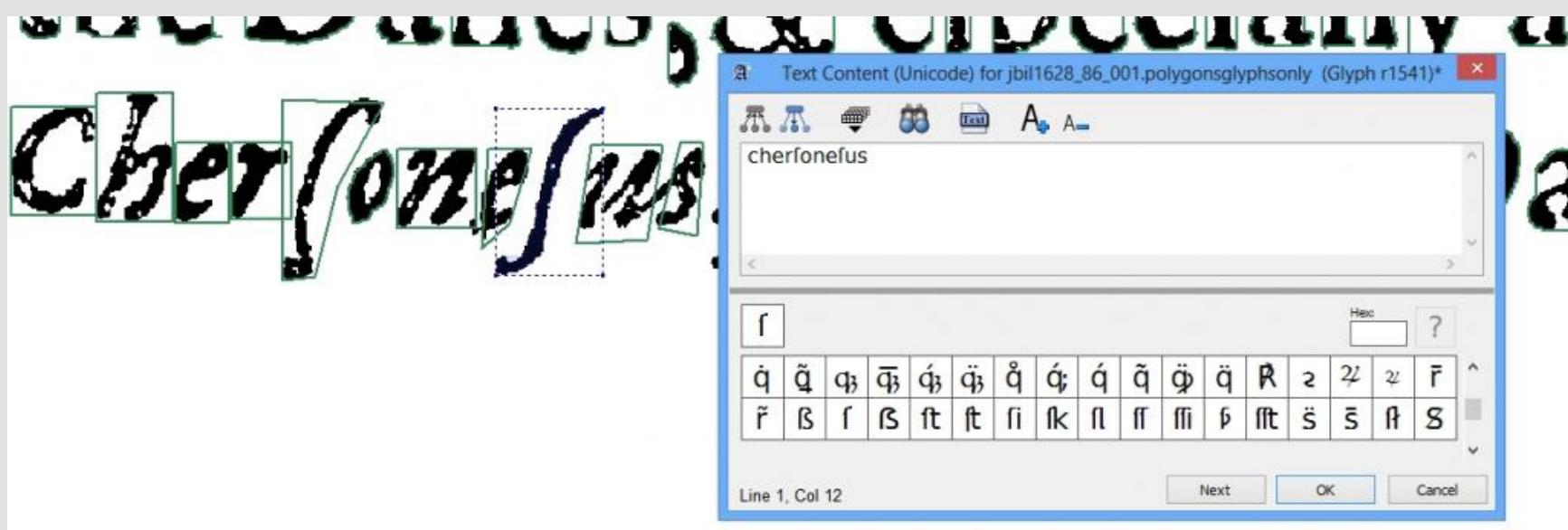
# Step 6: Text Correction

- There is difficulty reading letters that are distorted, obscured by noise or ink bleed-through, or if they are especially faint.



# Step 6: Text Correction

- Text misreads must be corrected by typing in the correct text in the Text Content Box for each corresponding character:



# Font Sets

To date the Aletheia team at the IDHMC has produced fourteen font training sets

Font Name	Publisher	Publication Year	Images Processed
JDAY1559_001	John Daye	1559	15
MFLE1633_94_017	Miles Flesher	1633	11
TCOT1632_82_001	Thomas Cotes	1632	11
IJAGG_guyot_83_99	François Guyot	Unknown	13
BRA_2788_002	Pierre Haultin	Unknown	8
cbil1692_110_002	Charles Bill	1692	6
cbil1693_116_001	Charles Bill	1693	10
cbil1702_116_1_001	Charles Bill	1702	14
hhil1686_82_1_001	Henry Hills	1686	15
jbill1628_86_001	John Bill	1628	15
IJAGG_guyot_83_15_B italics	François Guyot	Unknown	12
BRA_2788_004_B italics	Pierre Haultin	Unknown	5
MFLE1633_94_017_B italics	Miles Flesher	1633	11
TCOT1632_82_001_B italics	Thomas Cotes	1632	14

# Training Tesseract using Aletheia

- Once the font training libraries were in place, we started to train Tesseract by running each prepared font set through the Tesseract engine so that it would learn to read and recognize the characters of those fonts.

# Aletheia and Tesseract

## Aletheia Glyph Identifier

IR ffzftRcNf TRA-

My 12□ 800C1 RCCISWPF CO CUUCPUU 1232 \*y-

Whzc w2\$ 2w fitj 12 kz# 2U CMU Sj\*z#

OB I UxjNR yOU; OWOUM hcZ NUUC wxc GUWY

EUC OUCC kCfOz-C 1 FOICC CO ckb pUx-pzzf#p wkc\* 0

NU\*ICC Uzc hzUffc 2 1 10U8~

278

VVRC Lzfly

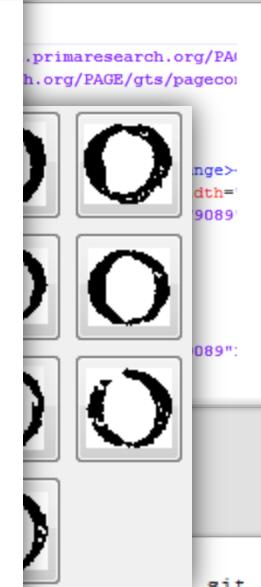
POB- N#MUXCUC-

ffNp N2y' kUC yOU w111?

## Tesseract Box

2	o	56	1844	82	1871	0
3	s	88	1840	106	1870	0
4	t	112	1844	127	1873	0
5	H	152	1844	193	1885	0
6	o	200	1844	226	1871	0
7	n	232	1844	257	1871	0
8	o	264	1844	290	1871	0
9	u	296	1844	320	1871	0
10	o	328	1844	350	1871	0

A more then common Stock ought to Prc  
 The Honours due to Smith's Immortal N  
 A Name whose early glorys were so hur  
 About ev'n in the Non-age of the Worl  
 That other Families were hardly knowr  
 Many a nobel Family had a nobler Name

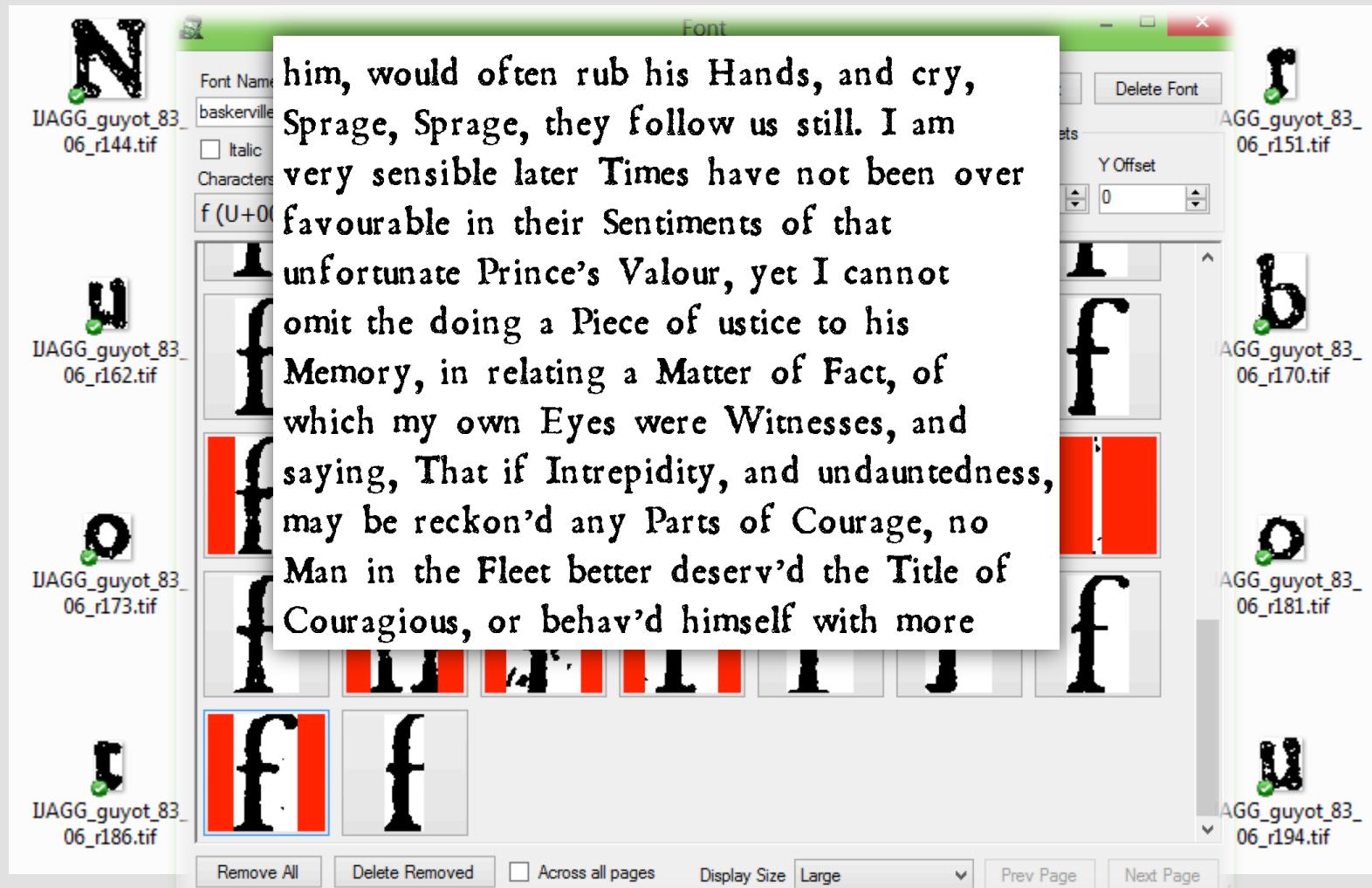


primaresearch.org/PAn  
 h.org/PAGE/gts/pageco:  
 nge>  
 dth= 9089  
 089":

# Roadblocks:

- How do we “feed” Tesseract only pristine instances of each letter for training purposes?
- How do we ensure that each glyph has been identified correctly and consistently?
- How do we eliminate the misidentification of pixel artifacts as characters?

# The Solution: Franken+



# Promising Results

## Before Franken+

OB I UxjNR yOU; OWOUM hcZ NUUC woxc GUWY  
 EUC OUCC kCfOz-C 1 FOICC CO ckP pUx-pzzfPp wkct 0  
 NU\*ICC Uzc hzUffc 2 1 10U8~

278  
 VVRC Lzfly TC hcc LOxch YOUPIC \*y 2  
 POB- N#MPUXCUC-  
 ffNP N2y' kUC yOU w111?

INff- I U22y UOC vCKCIy-  
 HCp Vc\$cly?

YOU FUC UU Offwjck IRkcx VOWCY hU\$ L-  
 ThQYUBh YOU WOUM fCCIC cPvffphcxc chc SC2\$\$ wlck OYkt  
 SkOUM ycc f2y' S\$UO 80IU8; VCzCIy

## After Franken+

That when thou meet'af& one, with enqtziring eyes  
 Do'i t search, and like a needy broker ptize

The Glke, and gold he weares, and to that rate  
 Sohigh or low, doG: raife thy formall hat:

That wilt confort none, Untill thou have knowne  
 What lands hee hath in hope, or of his owne,

As though all thy companions Ishould make thee  
 gointntes, and marry thy deare company.

Why HaouEd'R thou that dofk Uot onely apptovve;  
 But in ranke itchie luf?, deGre, andlove

The nakednefe and barrenneHe to enjoy,

# Thank You!

- Kathy Torabi: [torabik@tamu.edu](mailto:torabik@tamu.edu)
- Bryan Tarpley: [bptarpley@tamu.edu](mailto:bptarpley@tamu.edu)
- The IDHMC: <http://idhmc.tamu.edu>
- The EMOP Project: <http://emop.tamu.edu>