# EMOP

## EARLY MODERN OCR PROJECT

Start new document

TIFF Image

What's new:

• New tools for lines, words, glyphs
• Integrated OCR (Tesseract 3.02)
• ALTO and FineReader XML support
• Text overlay
• Drag and drop files into Aletheia

Reading Order and Layers

Validation

Open document

PAGE XML

Recent Documents:

cbil1692_110_002boxfilesglyphsonly.x

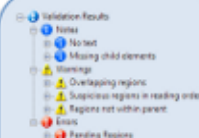BRA2788_002.boxfilesglyphsonly.xml

IJAGG_guyot_83_09polygonglyphsonly

IJAGG_guyot_83_09glyphsonly.xml
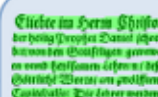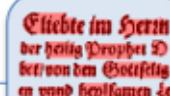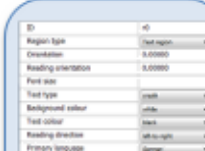
Document Image

Layout Regions
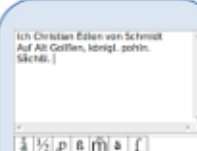
Text Lines

Words

Glyphs

B/W Image

Attributes

Text Content

# Step One: Naming Font Sets

Page images are named according to publisher, publication year, and the 20-line height. For example, a text image from the 1702 *Anno Regni text would be named cbil1702_116.*

| Printer | Primary Type Style | 20-LN HT (mm) | Body Size | Year | Place | Author | Title |
|---------|-------------------|---------------|-----------|------|-------|--------|-------|
| | | | | | | | |
| Andrew Crook | Roman | 143 | Double Pica | 1651 | London | Hobbes | Leviathan |
| | Roman | 92 | English | 1651 | London | Hobbes | Leviathan |
| | Italic | 92 | English | 1651 | London | Hobbes | Leviathan |
| Charles Bill | Black Letter | 116 | Great Primer | 1702 | London | England and Wales | Anno Regni |
| | Roman | 116 | Great Primer | 1702 | London | England and Wales | Anno Regni |
| | Black Letter | 116 | Great Primer | 1693 | London | England and Wales | Gulielmi et Mariae |
| | Roman | 116 | Great Primer | 1693 | London | England and Wales | Gulielmi et Mariae |
| | Roman | 94 | English | 1693 | London | England and Wales | Gulielmi et Mariae |
| | Black Letter | 116 | Great Primer | 1693 | London | England and Wales | Gulielmi et Mariae |

# Step Two: Page Segmentation and Analysis

Layout regions, lines, words, and individual glyphs in a text are identified and defined using Aletheia Desktop, a tool developed by PRImA Lab at the University of Salford that performs page segmentation and text recognition, among other tasks.

¶TO THE MOST REVE-
rend Father in God, and his sin-
giler good Lord, Edmond, by the
grace of God, Archbishop of Can-
terburie, Primate and Metropoli-
tane of all England, Thomas Timme
wisheth in grace and mercie of the
Spirite, in Christ Iesu.

Fter long exercise
in translating such
Latine Commenta-
ries vpon the holy
Scriptures, as I
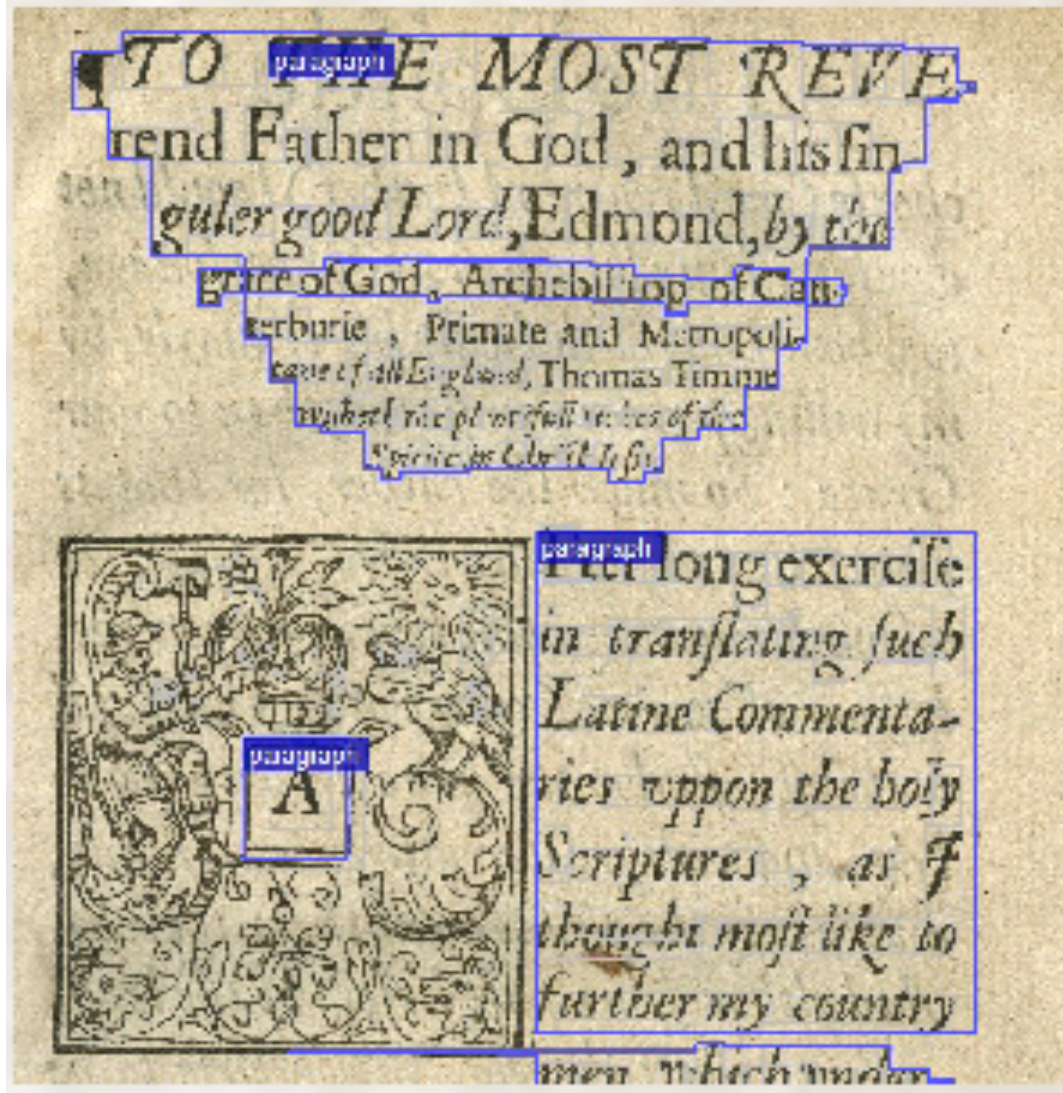thought most like to
further my country
men, which vnder-
stand not that tongue, to the said knowledge
of true Religion : at last I tooke in hand
M. Caluins exposition vpon Saint Pauls
Epistles to the Corinthians. And, as in my
poore iudgment, the writer is a most excel-
lent instrument of God, for the simple set-
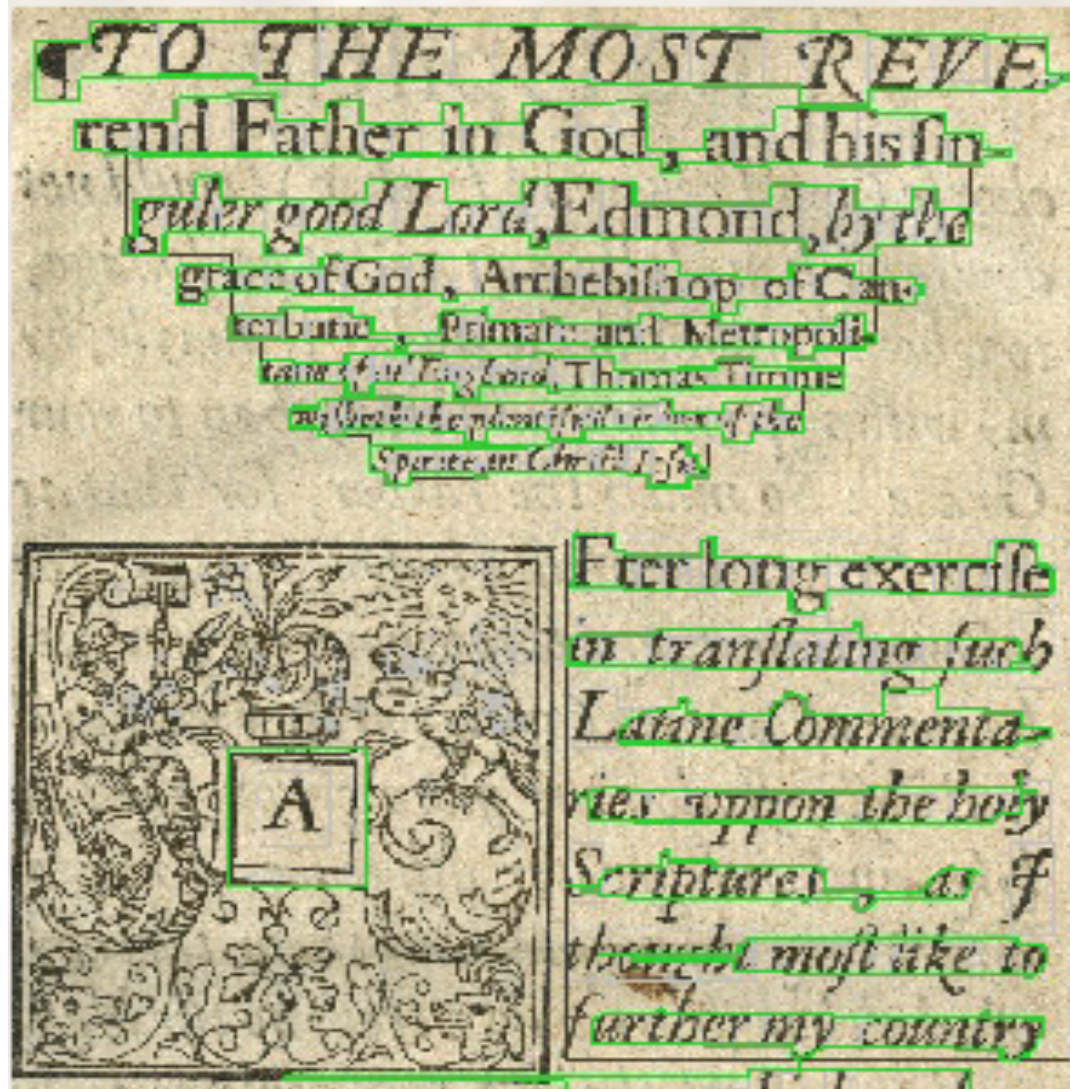ting foorth of his trueth, so in making my
¶. ij.                              choyse

# Layout Regions

# Text Lines

TO THE MOST REVE-
rend Father in God, and his sin-
guler good Lord, Edmond, by the
grace of God, Archebishop of Can-
terburie, Primate and Metropoli-
taine of all England, Thomas Timme
wisheth the gifts of the holy Ghost and the
Spirite in Christ Iesus.

Fter long exercise
in translating such
Latine Commenta-
ries vppon the holy
Scriptures, as I
thought most like to
further my country

# Words

# Glyphs

¶ TO THE MOST REVE-
rend Father in God, and His sin-
guler good Lord, Edmond, by the
grace of God, Archbishop of Can-
terburie, Primate and Metropoli-
tane of all England, Thomas Time me
wisheth the premisent riches of the
Spirite, in Christ Iesu.

A

Fter long exercise
in translating such
Latine Commenta-
ries vppon the holy
Scriptures, as I
thought most like to
further my country
men, which vnder-
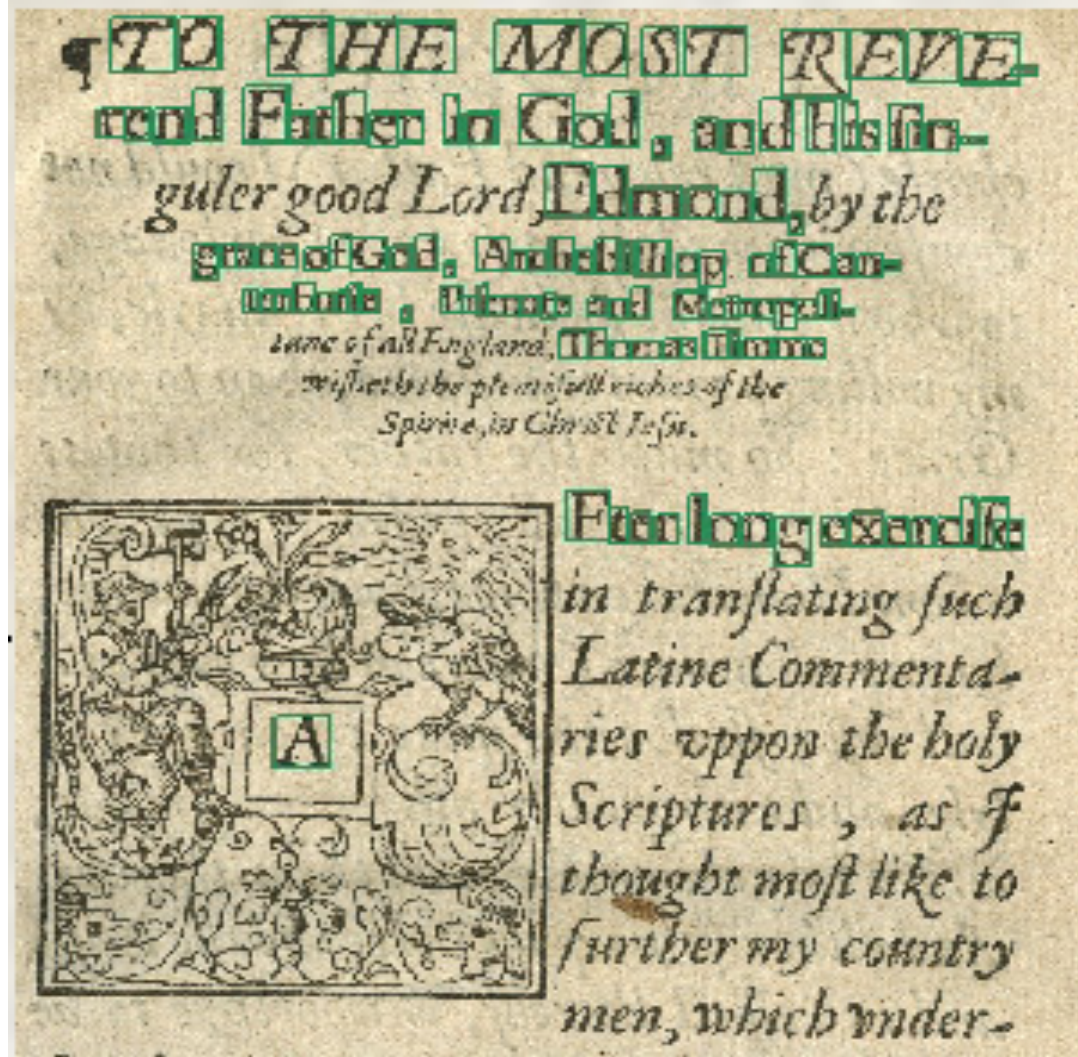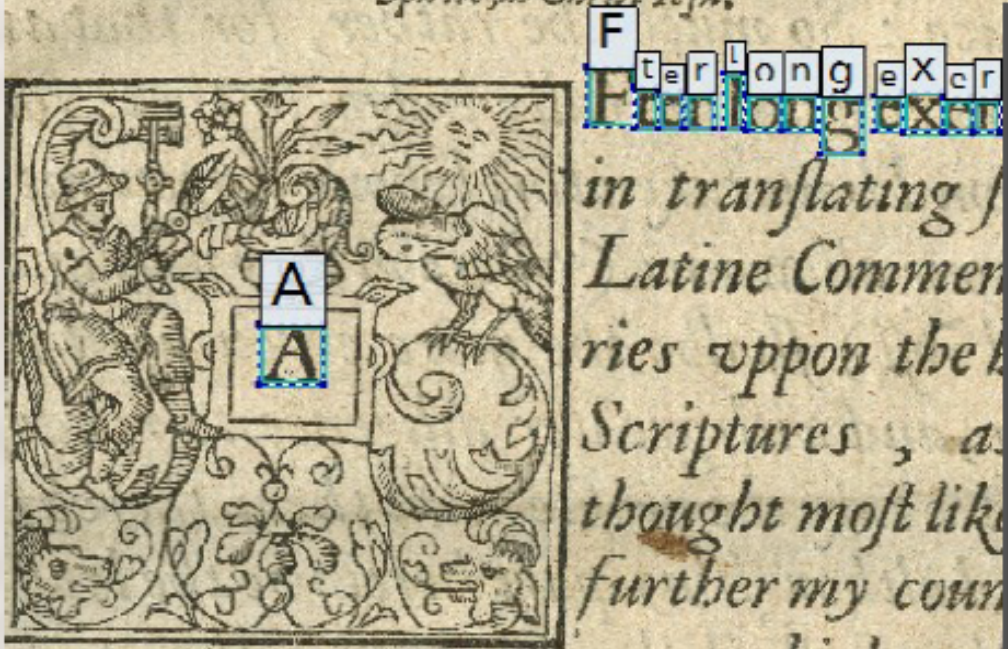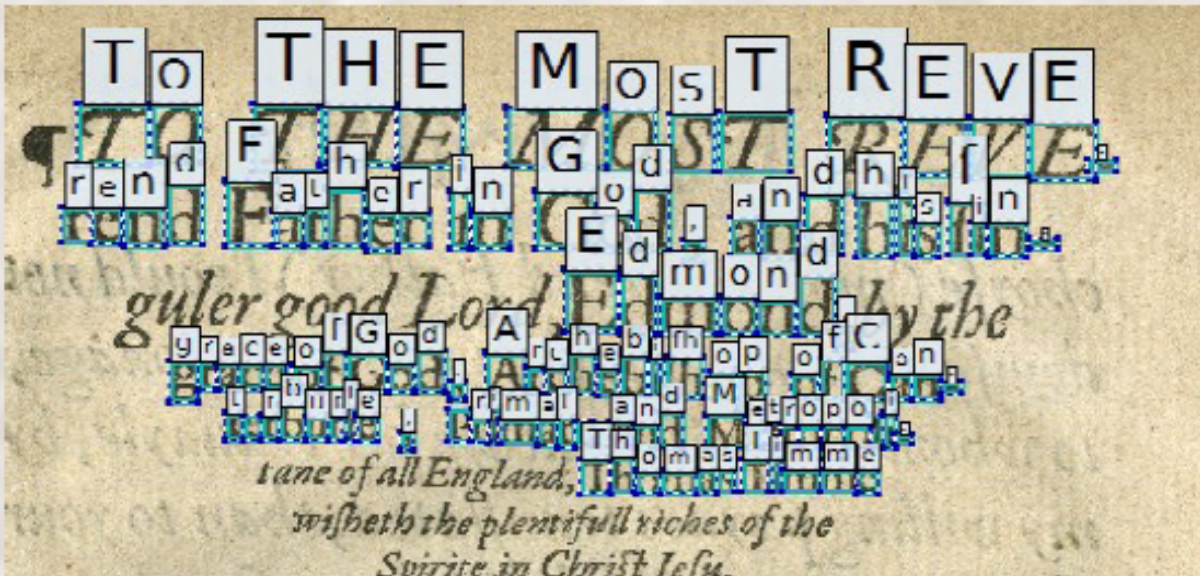
# Step Three: Text Correction

Aletheia analyzes the page image and generates text for each character. Students are then responsible for fixing any "misreads" by typing in the correct text in the Text Content Box for each corresponding character--paying special attention to unique characters such as long S's, ligatures, italics, rotunda R's, suspension marks, and printer's marks.
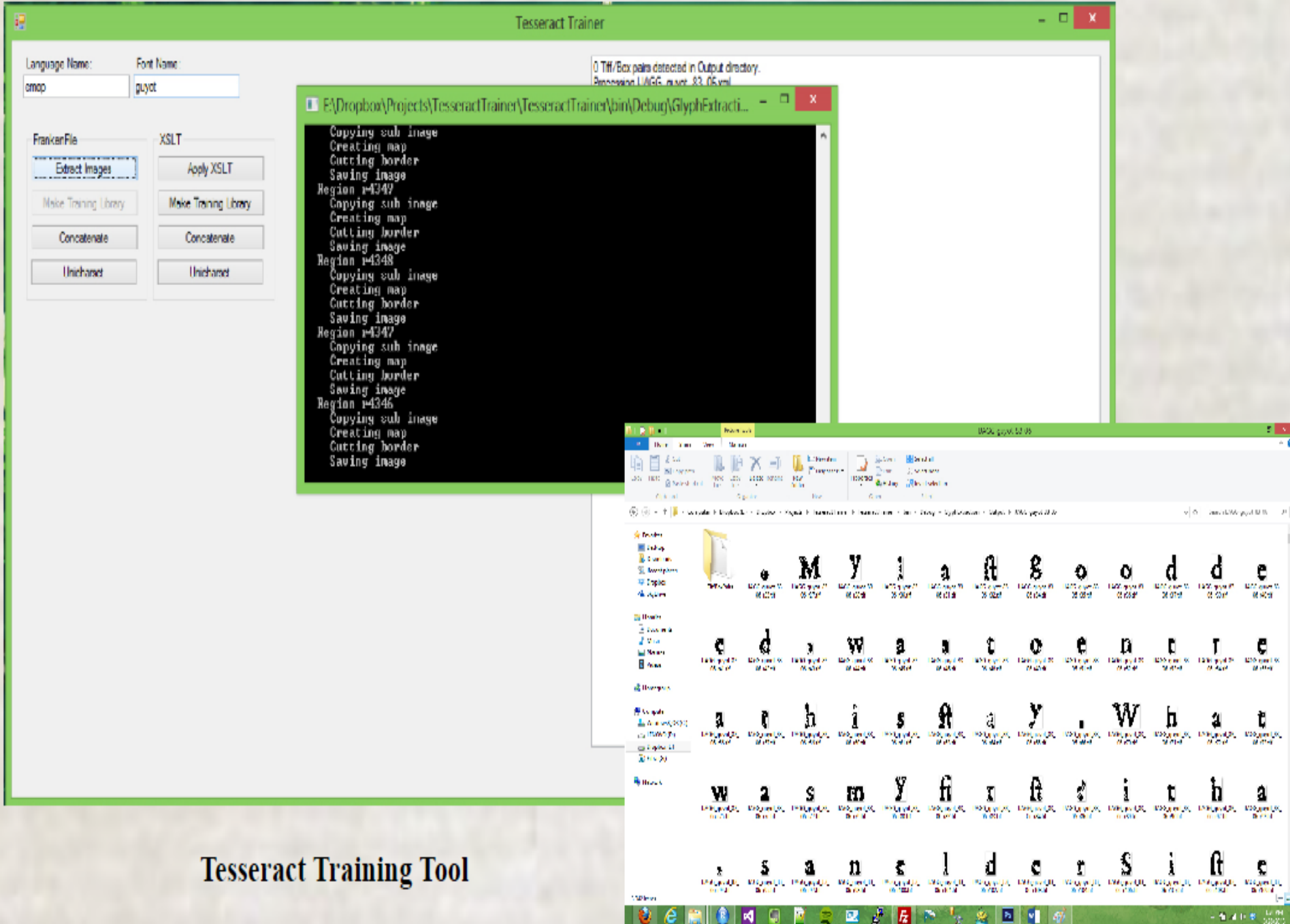
¶ TO THE MOST REVE-
rend Father in God, and his fin-
guler good Lord, Edmond by the
grace of God Archbishop of Can-
terbury, primat and Metropo-
litane of all England, Thomas Timme
wisheth the plentifull riches of the
Spirite, in Christ Iesu.

After long exer-
cise in translating
Latine Commenta-
ries vppon the holy
Scriptures, as I
thought most like
further my coun

**Text Correction**

To The Most Reverend

# Step Four: OCR Training

10-15 page images are processed for each font type and comprise a training set. A Tesseract training tool uses the font set to train the OCR engine to read and recognize that particular font.

**Tesseract Training Tool**