

DIGITIZATION OF HANDWRITTEN LEDGERS

AN INTEGRATED APPROACH TO IMPROVE HANDWRITTEN TEXT RECOGNITION

SUBMITTED IN PARTIAL FULFILLMENT FOR THE DEGREE OF MASTER OF SCIENCE

RODERICK MAJOOR
12852724

MASTER INFORMATION STUDIES
DATA SCIENCE
FACULTY OF SCIENCE
UNIVERSITY OF AMSTERDAM

SUBMITTED ON 23.05.2024

	UvA Supervisor
Title, Name	Dr. N.J.E. (Nanne) van Noord
Affiliation	MultiX Lab
Email	n.j.e.vannoord@uva.nl



1 **ABSTRACT**

2 Write your abstract here.

3 **KEYWORDS**

4 keywords, belong, here, with, commas, like, this

5 **GITHUB REPOSITORY**

6 <https://github.com/roderickmajoer/Thesis>

7 **1 INTRODUCTION**

8 The Bank of Amsterdam was an early bank established in 1609, de-
9 scribed by some as the first central bank [20]. The bank played an
10 important role in the financial world of the 17th and 18th century.
11 Transactions made at the bank from 1650 to 1800 were recorded in
12 ledgers and are still preserved to this day. These ledgers contain
13 information about transactions, both debit and credit, of customers
14 of the bank. The ledgers can thus be very important in analyzing
15 the money flow at the time. Digitizing these ledgers may help in
16 preserving them in a much easier to analyze format which would
17 allow for more research into the contents of the ledgers. However,
18 these ledgers are all handwritten making it hard to retrieve the
19 information out of them on large scale. Current efforts to retrieve
20 the information are done by using Handwritten Text Recognition
21 (HTR) tools to extract the text in the ledgers. The problem is that
22 the currently used HTR techniques have too many inaccuracies to
23 be used effectively. Because of this, the handwritten text cannot
24 be recognized correctly and can thus not be digitized properly. To
25 improve digitization efforts, it is crucial to find out why some hand-
26 written text cannot be recognized, what kind of errors are made
27 during the HTR process and how the HTR process can be optimized
28 to improve these inaccuracies. The research question we wish to
29 answer is:

31 *How can the categorization of errors in handwritten ledger analysis
32 be used to enhance the identification and resolution of text recognition
33 inaccuracies?*

34 Subquestions belonging to this research question are:

- 35
- What factors contribute to text recognition errors in hand-
written ledgers?
 - What are the different types of errors encountered in hand-
written text recognition?
 - How can specific processing methods improve text detection
and segmentation in handwritten ledgers?
 - What strategies can be implemented to mitigate common
text recognition errors in handwritten ledgers?

44 Since we perform a number of different experiments, our re-
45 port follows the following structure: We start by comparing some
46 related work. Then we will provide a description of the dataset.
47 After that we look at the different experiments, where for each
48 we provide our methodology including evaluation steps as well as

49 results and discussion of that particular experiment. We end with a
50 more general discussion and conclusion.

51 **2 RELATED WORK**

52 The research gap being addressed in this project revolves around the
53 challenges associated with the digitization of handwritten ledgers.
54 The proposed research will contribute in closing the research gap
55 within the HTR domain, specifically looking at historical handwrit-
56 ten documents. The digitization of handwritten ledgers contains
57 several steps. In this section, we will look at key research papers
58 that discuss different approaches to the steps of our problem.

59 **2.1 Document Layout Analysis**

60 The biggest challenge in our digitization task is being able to find the
61 correct word order from the HTR output. Seuret describes layout
62 analysis and finding the correct word order as one of the main
63 challenges in HTR [21]. Several approaches to solve this problem
64 have been tried. These approaches are typically divided in either
65 deep learning methods or more classical computer vision methods.
66 We will show recent advancements made in both approaches.

67 Kastanas et al. show the use of different deep learning archi-
68 tectures for layout analysis [8]. Transformer-based, graph-based
69 models, and convolutional neural networks (CNN) are compared.
70 The CNN-based model YOLOv5 and transformed-based model Lay-
71 outLMv3 both show promising results for identifying most elements
72 on documents, with YOLOv5 performing slightly better [8].

73 Multiple other studies describe the use of transformer based
74 methods to perform layout analysis and optical character recog-
75 nition (OCR) [10, 13, 24]. Furthermore, Smock et al. show that
76 transformer-based object detection models can produce excellent
77 results for the tasks of detection, structure recognition, and functional
78 analysis of tables [22]. While these approaches show promising
79 results, they are not tested for our specific task at hand.

80 Oliveira et al. propose dhSegment, an open-source implemen-
81 tation of a CNN-based pixel-wise predictor for document segmen-
82 tation. The approach aims to address various document processing
83 tasks simultaneously, including page extraction, baseline extraction,
84 layout analysis, and illustration and photograph extraction. They
85 show that a single CNN architecture can be used across tasks with
86 competitive results [2].

87 Drobny et al. propose a holistic method that applies Mask R-CNN
88 for text line extraction in historical documents. Their work achieved
89 state-of-the-art results on well known historical datasets [6]. While
90 their work does not directly transfer to our task, it does show the
91 potential for document segmentation and possibly table recogni-
92 tion abilities of Mask R-CNN.

93 The aforementioned studies show that the use of deep learning
94 methods can be used to perform layout analysis as well as table
95 recognition tasks in historical documents. However, these methods
96 often require massive amounts of training data to be used effectively.
97 To use these methods, we should look whether pre-trained models
98 transfer well to our dataset and are able to detect the layout of our

99 documents. Other approaches use more classical computer vision
100 techniques to segment a document in order to recognize the layout.
101 These methods often require homogeneity of the documents to
102 work effectively.

103 Lehenmeier et al. show a method that combines various state-
104 of-the-art methods for OCR pro- cessing to detect layout on a
105 document [9]. To perform table recognition, they make use of tech-
106 niques such as binarization, denoising, and Hough line detection
107 to find relevant parts of the tables. By taking the intersection of
108 relevant horizontal and vertical Hough lines they are able to find
109 table cells and thus determine the layout of the table [9].

110 Bulacu et al. introduce a method based on contour tracing that
111 generates curvilinear separation paths between text lines in order to
112 preserve the ascenders and descenders of text lines to determine the
113 layout in handwritten historical documents [3]. With this approach
114 they are able to determine seperate elements of tables in historical
115 documents.

116 Liu et al. describe the use of several key steps, including skew
117 correction, region segmentation, and page layout analysis on his-
118 torical Tibetan documents [12]. Skew correction is achieved by
119 leveraging baseline features and Hough transform to determine
120 the document's skew angle. Subsequently, region segmentation is
121 accomplished by identifying borders through a series of preprocess-
122 ing steps, including median filtering, Gaussian smoothing, Sobel
123 edge detection, and removal of small area regions. These processes
124 collectively enable accurate positioning of document borders for
125 further analysis and structure extraction [12].

126 Liang et al. describe the use of Hessian and Gabor filters in
127 order to find table structures on a document [11]. They show the
128 possibility for column segmentation for the tables using the found
129 lines.

130 Prieto et al. address the challenge of document image understand-
131 ing in documents with complex layouts like tables. They compare
132 two approaches and show promising results [19]. Their approaches
133 take the text lines outputted by HTR systems and use machine
134 learning methods to group these text lines and classify cells based
135 on these grouped text lines in order to find the tabular structure.

136 It can be seen that depending on the dataset and the homogene-
137 ity of the documents, different approaches for layout analysis/table
138 recognition may be suitable. When a similar table structure is clearly
139 drawn on all the documents, approaches using basic image process-
140 ing techniques could be suitable to determine the structure of the
141 table. When the structure is not easily identifiable on the page, it
142 might be better to look at approaches that use the output of HTR
143 systems such as text lines or words locations to classify table cells.
144 Machine learning based approaches could also be useful in these
145 cases but often require many training data.

2.2 HTR Processing Techniques

146 Typical HTR pipelines consist of pre-processing an image to make
147 it ready for HTR, running the image through the HTR system and
148 post-processing the output of the HTR system [9].

149 Studies have shown that pre-processing images can lead to better
150 results when using them as input for document layout analysis or

152 OCR. Thus, we want to optimize our image material such that the
153 document layout analysis and HTR can be performed more accu-
154 rately. There are different methods for doing this such as adjusting
155 the contrast to remove artifacts and noise or fixing the image align-
156 ment [1, 5, 9]. It is crucial to find out which pre-processing steps
157 might be useful to our problem at hand. Chen et al. show that the
158 use of several pre-processing steps including character segmenta-
159 tion, tilt correction, offset correction, size normalization and image
160 thinning lead to better HTR output results [4].

161 Post-processing OCR and HTR output typically consists of trying
162 to correct misidentified words or characters. One strategy to do
163 this is to create a dictionary of candidate characters and use the
164 Levenshtein distance to calculate the distance between predicted
165 and dictionary characters [18]. The use of a dictionary for allowed
166 characters is shown to be a good strategy to eliminate HTR errors
167 [9]. Other studies show the use of language models to correct word
168 output and spelling in HTR systems [15, 18].

2.3 HTR Evaluation Strategies

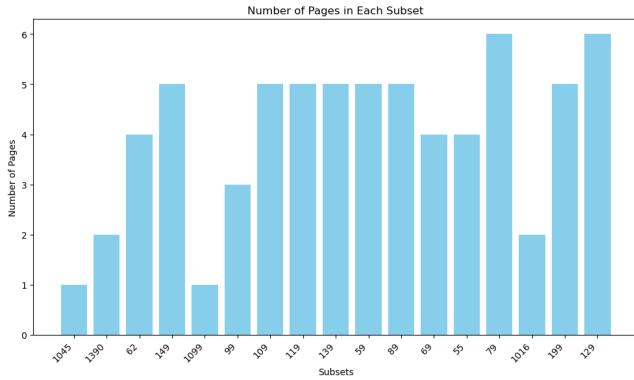
169 Typically, HTR systems are evaluated using the word error rate
170 (WER) and character error rate (CER). However, the evaluation
171 of page-level HTR output faces challenges primarily due to the
172 Reading Order (RO) problem [23]. In their study, they define a bag-
173 of-words Word Error Rate (bWER) to accurately measure page-level
174 RO-independent word errors as well as a regularized version of the
175 Hungarian Algorithm (HA) to compute word- and character-level
176 RO-independent recognition accuracy [23]. Another strategy is
177 the use of the intersection over union (IoU) score which matches
178 ground truth and HTR output boxes based on their coordinates
179 on the page [17]. It is then possible to use the WER and CER to
180 calculate the accuracy.

3 DATA DESCRIPTION

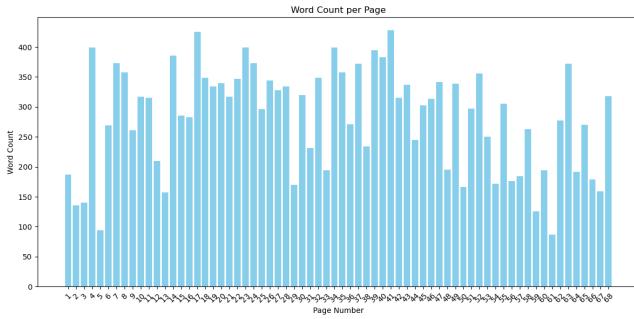
182 The dataset used in this study consists of pages from the collection
183 of ledgers from the Bank of Amsterdam. The full collection can be
184 found on the [website](#) of The Amsterdam City Archives. In the series
185 of ledgers, the current accounts of the customers were maintained.
186 Each customer had one or more pages, with smaller traders having
187 a portion of a page. An alphabetical list of the traders and the page
188 of their current account can be found in the index. The ledgers kept
189 track of the amounts that were debited and credited from and to a
190 customer. The ledger pages consist of several columns containing
191 information such as the date, account holder name, account number
192 and amount debited/credited. Figure 9 shows an example of what
193 a page looks like. Our dataset is split into multiple subsets that
194 come from different books and that all have slightly varying pages,
195 while still having structural similarity. Figure 1 shows the amount
196 of pages for each subset in the dataset. The total dataset consists of
197 68 pages. A word count for each page is shown in Figure 2

199 The ground truth and HTR system output both use a [PageXML](#)
200 format to store values on the page. These PageXML store the coor-
201 dinates for each item (e.g. textregion, tableregion, etc.) on the page.
202 We can plot the values in the PageXML file over our image to better
203 show this. Figure 10 shows the PageXML data of the ground truth

204 and HTR system plotted on the original image for part of one page
205 in our dataset.



241 **Figure 1: The amount of pages for each subset in our dataset.**



241 **Figure 2: The word count for each page in our dataset.**

226 and HTR output based on the highest IoU score, we can compare the
227 ground truth values and HTR output on a word to word basis.

228 Once we have aligned the words in the ground truth and the HTR
229 output properly, we can start to identify the flaws of the HTR system.
230 We will make use of the Levenshtein distance, which can be used
231 to find the minimum number of single-character edits (insertions,
232 deletions, or substitutions) required to change one string into the
233 other [7]. By applying this on our ground truth values and HTR
234 output on a word to word basis, we can track how many character
235 insertions, deletions and substitutions are done in the HTR output
236 for each word in the ground truth. We will show this for each subset
237 of our data to also see whether some subsets are more difficult for
238 the HTR system. We will also keep track of what characters are
239 inserted, deleted or substituted so we can find what characters seem
240 to be most problematic for the HTR system.

241 By creating this analysis, we can also better inspect our image
242 by visually showing the words that the HTR system can or cannot
243 match. An example of this visual representation is shown in Figure
244 11. Using this visual representation, we can inspect the image to
245 see whether there are any characteristics that might cause the HTR
246 system to make errors. This could be noise on parts of the page,
247 words being close together or even overlapping etc.

248 Upon completing our analysis, we can conduct a list of main
249 errors made by the HTR system and decide on what steps could
250 be taken to mitigate these errors. The approach for these possible
251 solutions is outlined in further sections.

252 **4.2 Evaluation**

253 For the evaluation of the HTR system, metrics such as WER, CER,
254 and Levenshtein Distance will be calculated. These metrics measure
255 the minimum number of operations (substitutions, insertions,
256 deletions) needed to transform the HTR output into the ground
257 truth. The formulas for these metrics are:

$$258 \quad WER = \frac{S_w + D_w + I_w}{N_w} \quad (1)$$

259 Where S_w , D_w , I_w are the number of word substitutions, deletions
260 and insertions respectively and N_w is the total number of words in
the ground truth.

$$261 \quad CER = \frac{S_c + D_c + I_c}{N_c} \quad (2)$$

262 Where S_c , D_c , I_c are the number of character substitutions, deletions
263 and insertions respectively and N_c is the total number of characters
in the ground truth.

264 We will also calculate recall and precision scores. To calculate
265 those for HTR output, we first need to define what constitutes true
266 positives (TP), false positives (FP), and false negatives (FN) in our
267 context.

- 268 • **TP (True Positive):** Characters correctly recognized by the
269 HTR system.
- 270 • **FP (False Positive):** Characters incorrectly recognized by
the HTR system (i.e., system identifies text where there is
none or incorrectly identifies text).

- 273 • **FN (False Negative):** Characters that are present in the
 274 ground truth but not recognized by the HTR system.

275 The recall and precision are then:

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

$$Precision = \frac{TP}{TP + FP} \quad (4)$$

276 4.3 Results

277 The accuracy of the current HTR system using the aforementioned
 278 metrics is shown in Table 1. A more in-depth analysis about the
 279 amount of character insertions, deletions, substitutions and matches
 280 made by the HTR system can be seen in Figure 5 and Figure 3.
 281 Figure 4 shows the frequencies of characters inserted, substituted
 282 and deleted. We use these analyses to conduct a list of errors made
 283 by the HTR system. A selection of the most common errors and
 284 possible solutions is shown in Table 2.

Table 1: Baseline scores for the currently used HTR system.
 We matched the ground truth words with HTR output words
 based on highest IoU value.

Metric	Value
WER	0.4243
CER	0.3422
Recall	0.9475
Precision	0.7403
Total Edit Distance	18920

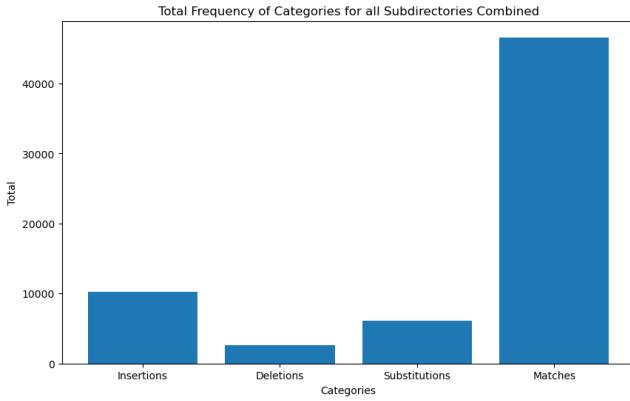
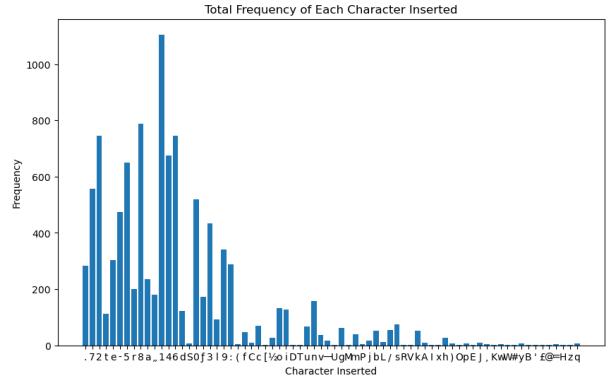
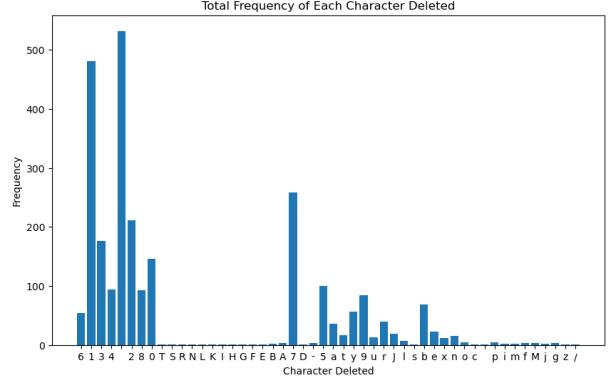


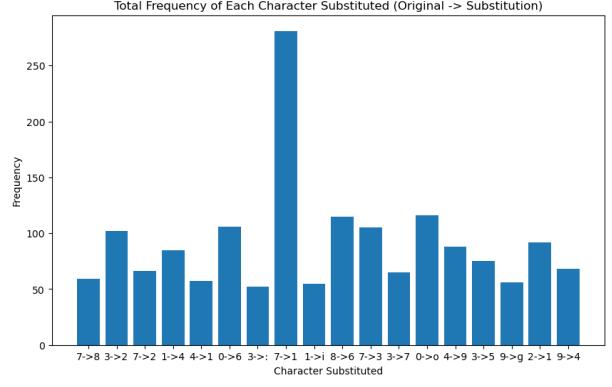
Figure 3: The total amount of insertions, deletions, substitutions and matches made by the HTR system compared to the ground truth values. The values are on a character level and made word per word.



(a) Frequency of inserted characters made by the HTR system.



(b) Frequency of deleted characters made by the HTR system.

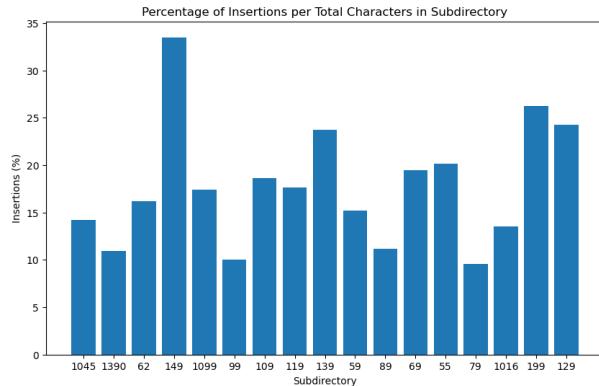


(c) Frequency of substituted characters made by the HTR system.
 The substitutions are shown as (original character -> substituted character). We only show the substitutions with a frequency of 50 or higher.

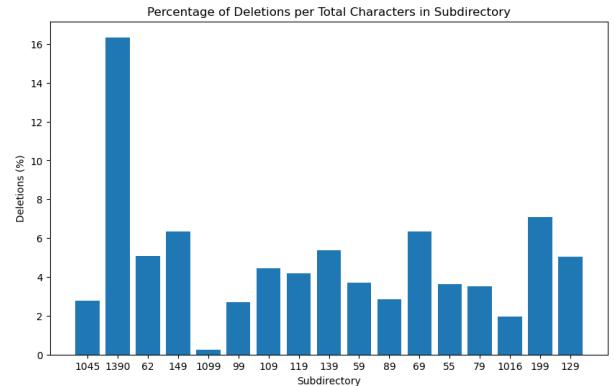
Figure 4: The frequencies of inserted, deleted and substituted characters made by the HTR system.

Table 2: Categorization of HTR System Errors

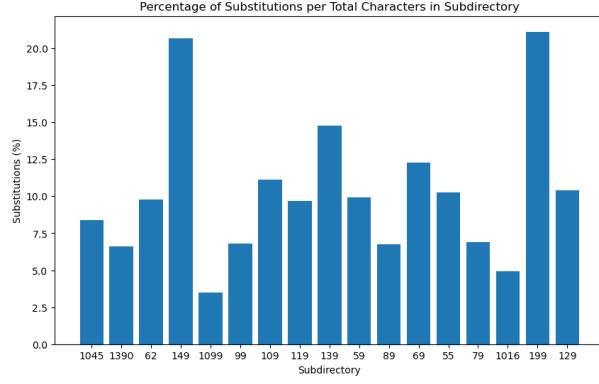
Error Type	Description	Possible Solution
Deletions	Characters not recognized at all (false negative)	Pre-process image, Improve HTR system
Insertions	Characters recognized that are not there (false positives)	Pre-process image, Improve HTR system
Substitutions	Characters recognized as another character	Pre-process image, Improve HTR system
Layout - Word Split	Single word is split into separate words	Layout analysis
Layout - Word Merge	Multiple words are recognized as one word	Layout analysis
Layout - Data Split	Data wrongly split (e.g., 'feb 5' in GT becomes 'feb' and '5' in HTR)	Layout analysis
Layout - Word Order	Word order not recognized	Layout analysis
HTR - Bounding Box	Wrong location bounding box vs actual location word	Improve HTR system
Case Sensitivity	Case sensitivity issues	Post-process string
Invalid Characters	Characters that cannot occur	Post-process string



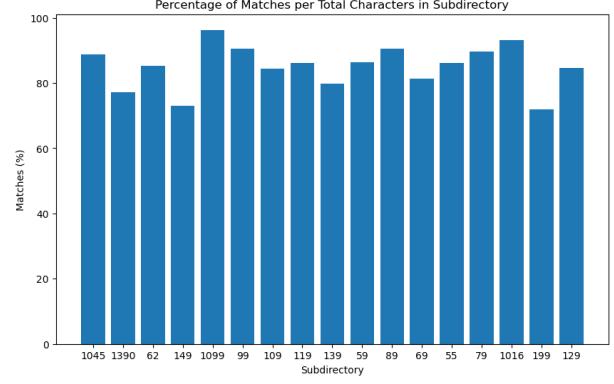
(a) Amount of insertions made by the HTR system per subdirectory.



(b) Amount of deletions made by the HTR system per subdirectory.



(c) Amount of substitutions made by the HTR system per subdirectory.



(d) Amount of matches made by the HTR system per subdirectory.

Figure 5: The amount of insertions, deletions, substitutions and matches made by the HTR system compared to the ground truth values. The values are on a character level and made word per word and as a percentage of total characters in ground truth.

285 5 POST-PROCESSING

286 5.1 Method

287 From our error analysis, we find that using post-processing tech-
 288 niques may lead to improved accuracy. Post-processing in our ex-
 289 perimental framework focuses on refining the output generated
 290 by the HTR system. One such post-processing step involves the
 291 removal of characters from the output strings that do not align with
 292 the ground truth dictionary [9, 18]. Analysis of the HTR output (as
 293 depicted in Figure 4a) reveals several insertions that are not valid
 294 letters or numeric characters. To address this, we employ regular
 295 expressions (regex) to filter out invalid characters from the HTR
 296 output. We are then left with alphanumeric characters only. Fur-
 297 thermore, the HTR system often substitutes letters for numbers,
 298 as can be seen in Figure 4c. We can undo these substitutions in
 299 some cases. For instance, when a single character is found to be a
 300 letter by the HTR system, we can expect this to be wrong, since
 301 single character words are never letters in the ground truth. Also,
 302 numbers that contain a letter (such as '11i0') are probably wrong.
 303 We use these facts and the most common made substitutions as
 304 found in our earlier analysis to recover the right characters. We use
 305 the same evaluation metrics as before to compare our method to
 306 the baseline scores.

307 5.2 Results

308 We compare the post-processed HTR output to the original HTR
 309 output to see the different accuracy scores. The accuracy scores
 310 for both the baseline and post-processed HTR output are shown in
 311 Table 3. The comparison of the performances is shown in Figure
 312 6. We see a very slight decrease in recall while all other metrics
 313 improve.

Table 3: WER, CER, precision and recall scores for the HTR system output with and without post-processing techniques applied.

Metric	Baseline	Post-Process
WER	0.4243	0.3916
CER	0.3422	0.3169
Recall	0.9475	0.9418
Precision	0.7403	0.7622
Total Edit Distance	18920	17508

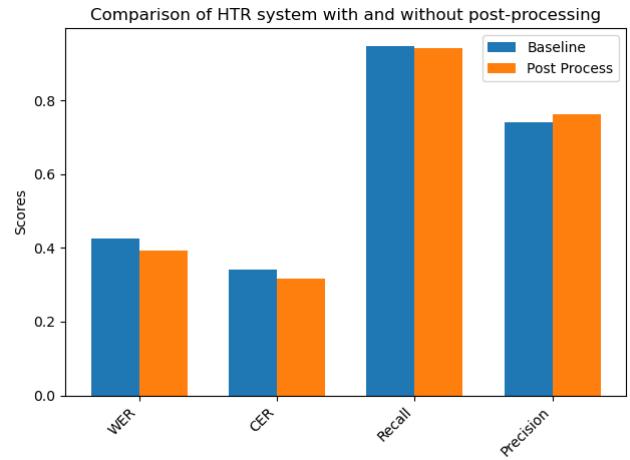


Figure 6: The WER, CER, precision, and recall scores for the HTR system with and without post-processing applied.

314 6 LAYOUT ANALYSIS

315 6.1 Method

316 Up until this point, we have been comparing the HTR output to the
 317 ground truth values by matching their bounding boxes based on the
 318 highest IoU score. We did this because the HTR system currently
 319 is not able to capture the layout of the the tabular structure of
 320 the ledger pages and thus we could not just match the output
 321 directly, since the word order would be wrong. The final part of
 322 the methodology involves analyzing the layout of the documents.
 323 Recognizing the layout is a crucial step in digitizing the ledger pages,
 324 since without it there is no good way to use the output of the HTR
 325 system because the correct word order would not be known. We will
 326 perform the layout analysis by using image processing techniques
 327 as well as making use of the HTR output. For this, we will use the
 328 [OpenCV \(cv2\)](#) library to perform the necessary processing steps.

329 To detect the tabular structure on the ledger pages, we must find
 330 a way to detect the columns and rows of the table to get the table
 331 cells. The ledger pages contain physical lines, which separate the
 332 different columns. We start by trying to detect these lines.

333 The first step is to pre-process our image. We do this by convert-
 334 ing our image to grayscale and applying a Gaussian blur on the
 335 grayscale image. This helps reducing the computational complexity
 336 and reduces high frequency noise [14].

337 Next, we apply two types of thresholding to the blurred image:
 338 Otsu's thresholding and Adaptive Gaussian thresholding. Otsu's
 339 thresholding only captures the most clear markings on the page,
 340 such as the text, while the more thin and faint column lines are not
 341 captured [3]. The Adaptive Gaussian threshold captures both the
 342 text as well as the column lines. By subtracting the first threshold
 343 from the latter, we are left with the column lines and some noise.

344 Following thresholding, we perform morphological operations,
 345 specifically dilation and erosion [16]. We use specific kernels to
 346 remove noise while retaining the vertical lines.

347 We then apply the Hough Line Transform to detect the lines
 348 in the image. This can be used to detect (nearly) straight lines in
 349 images [9]. We filter the detected lines based on their length and
 350 position. For each detected line, we calculate its slope and intercept,
 351 and then draw an extended line from the top to the bottom of the
 352 foreground area (which contains the ledger page) in the image.
 353 This is done using basic principles of line equations in coordinate
 354 geometry.

355 Finally, we perform additional dilation and erosion to merge
 356 close lines and remove final noise. We then find contours in the
 357 image and draw a line from the top to the bottom of each contour.
 358 We are then left with an image containing the contours of each
 359 column section, which can then be used to create bounding boxes
 360 for the different columns. The whole process is shown in Figure 12.

361 To detect the rows present in the ledger pages, we have to use
 362 a different strategy, since there are no physical markings for row
 363 lines on the pages. We will use the word bounding boxes which are
 364 outputted by the HTR system in order to find the row bounding
 365 boxes.

366 We start by extracting the bounding boxes of the words outputted
 367 by the HTR system. We will only use the word bounding boxes that
 368 are kept after performing the post-processing. Since the images in
 369 our dataset typically consist of two tables (debit and credit) we have
 370 to determine where both tables start and end. We do this by finding
 371 the middle column line from the list of column lines we found by
 372 our method above. We then separate the bounding boxes into two
 373 groups: left boxes and right boxes, based on their x-coordinates
 374 relative to the middle column line.

375 We then take a list of boxes (either left boxes or right boxes)
 376 and group them into rows. This is done by sorting the boxes by
 377 their x-coordinates (from left to right) and then grouping boxes
 378 with similar y-coordinates together. For each box, we find the best
 379 row to append the box to, based on the minimum difference in
 380 y-coordinates. If no suitable row is found, we create a new row. The
 381 result of these steps are shown in Figure 13.

382 Once we have found the rows, we can extend the bounding boxes
 383 from the left side of the ledger page to the middle for the left boxes
 384 and from the middle to the right side of the page for the right boxes.
 385 We are now left with both the column and row bounding boxes.

386 Once the row and column bounding boxes are determined, we
 387 can use the intersection of these bounding boxes to determine the
 388 table cells. Figure 7 shows an example of these found table cells.

389 6.2 Evaluation

390 For the layout analysis methods, we take the ground truth table cell
 391 bounding boxes and compare them to our method's predicted table
 392 cell bounding boxes. We do this by matching each ground truth cell
 393 to a predicted cell based on the highest IoU score between the two.
 394 Only one ground truth cell can be matched to one predicted cell
 395 and vice versa. We then determine whether cells are a true positive
 396 by setting an IoU threshold. If ground truth and predicted cells are
 397 matched with an IoU score higher than this threshold, we consider
 398 it a true positive. Ground truth cells that are not matched with an
 399 IoU score above the threshold are considered a false negative and

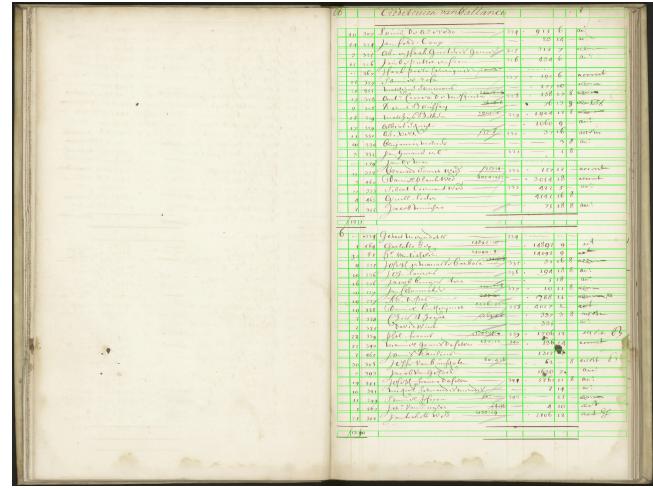


Figure 7: The bounding boxes of the predicted table cells after taking the intersection of the rows and columns.

400 predicted cells that are not matched with an IoU score above the
 401 threshold are considered a false positive. We will use a threshold
 402 of 0.5, 0.25 and >0. We can then determine the recall and precision
 403 scores, using the same formula as before. Additionally, we calculate
 404 the F1-score:

$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (5)$$

405 Furthermore, after matching the cells as described above, we
 406 further verify the correctness of the match by comparing the text
 407 content of the predicted table cells against the text content in the
 408 ground truth table cells. This is done by taking the words from the
 409 HTR output that match location with the predicted cells. The text
 410 comparison is done using a normalized Levenshtein distance with
 411 various thresholds: 0 (exact match), 0.25, 0.5, and 0.75. We use
 412 these thresholds because we cannot be sure that words from the
 413 HTR output match exactly with ground truth words because of the
 414 imperfections of the HTR system. For each threshold, we calculate
 415 the mean accuracy of the word matches. Additionally, we calculate
 416 the CER for all the words in the ground truth cells and the words
 417 in the predicted cells. If the cells are indeed matched correctly,
 418 we would expect to see good mean accuracy of the word matches
 419 (especially for the lower thresholds) and we would expect the CER
 420 to be close to the value found in Section 4.

421 6.3 Results

422 The performances using IoU thresholds of 0.5, 0.25 and >0 are
 423 shown in Table 4, Table 5 and Table 6 respectively. We show the
 424 performance both at the subset level as well as the total mean scores.
 425 The 'IoU Score' indicates the mean score for the true positives found
 426 at this IoU threshold. The comparison of the total mean scores can
 427 be seen in Figure 8. The accuracy of comparison of text contents
 428 for ground truth and predicted cells can be seen in Table 7.

Table 4: Table Cell Performance Metrics at IoU Threshold 0.5.

Subdir	IoU Score	Precision	Recall	F1
1045	0.638	0.674	0.674	0.674
1390	0.641	0.303	0.468	0.363
62	0.627	0.437	0.487	0.460
149	0.608	0.310	0.383	0.342
1099	0.621	0.532	0.499	0.515
99	0.638	0.525	0.504	0.514
109	0.645	0.520	0.552	0.534
119	0.617	0.422	0.437	0.429
139	0.626	0.435	0.475	0.453
59	0.664	0.555	0.541	0.548
89	0.632	0.466	0.519	0.490
69	0.626	0.391	0.456	0.420
55	0.631	0.287	0.319	0.302
79	0.614	0.389	0.439	0.412
1016	0.661	0.637	0.635	0.636
199	0.594	0.243	0.339	0.278
129	0.617	0.342	0.404	0.369
Mean	0.627	0.416	0.459	0.434

Table 5: Table Cell Performance Metrics at IoU Threshold 0.25.

Subdir	IoU Score	Precision	Recall	F1
1045	0.546	0.990	0.990	0.990
1390	0.498	0.513	0.785	0.613
62	0.471	0.854	0.959	0.902
149	0.440	0.710	0.879	0.784
1099	0.473	0.939	0.881	0.909
99	0.504	0.916	0.878	0.896
109	0.508	0.848	0.894	0.867
119	0.460	0.869	0.901	0.885
139	0.474	0.819	0.894	0.853
59	0.525	0.910	0.888	0.899
89	0.475	0.862	0.961	0.907
69	0.464	0.821	0.963	0.884
55	0.451	0.726	0.811	0.765
79	0.452	0.829	0.929	0.875
1016	0.553	0.835	0.831	0.832
199	0.432	0.639	0.894	0.731
129	0.453	0.738	0.878	0.801
Mean	0.473	0.804	0.898	0.844

Table 6: Table Cell Performance Metrics at IoU Threshold >0

Subdir	IoU Score	Precision	Recall	F1
1045	0.196	1.000	1.000	1.000
1390	0.152	0.603	0.952	0.729
62	0.190	0.876	0.984	0.925
149	0.169	0.762	0.943	0.841
1099	0.206	0.945	0.886	0.915
99	0.177	0.960	0.920	0.939
109	0.188	0.885	0.933	0.905
119	0.178	0.906	0.939	0.922
139	0.185	0.860	0.938	0.896
59	0.192	0.946	0.923	0.934
89	0.186	0.873	0.974	0.919
69	0.195	0.839	0.984	0.903
55	0.172	0.820	0.911	0.862
79	0.182	0.863	0.967	0.911
1016	0.232	0.845	0.840	0.841
199	0.169	0.690	0.962	0.788
129	0.164	0.802	0.952	0.869
Mean	0.182	0.845	0.946	0.887

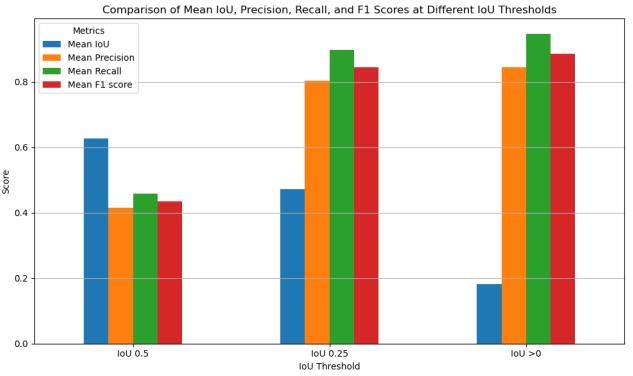


Figure 8: The mean IoU, precision, recall and F1 scores for ground truth and predicted table cell bounding boxes for different IoU thresholds.

Table 7: Mean accuracy and CER for matching ground truth table cells with predicted table cells using different IoU thresholds and comparing their text contents using normalized Levenshtein distance thresholds. Values for LT 0, 0.25, 0.50, and 0.75 indicate the mean accuracy at which the normalized Levenshtein distance between the ground truth and predicted texts is within 0 (total match), 25%, 50%, and 75% respectively. CER indicates the average character error rate for the text content in the matched cells.

IoU Threshold	LT 0	LT 0.25	LT 0.50	LT 0.75	CER
>0	0.5369	0.5996	0.7333	0.8242	0.4290
0.25	0.5555	0.6213	0.7563	0.8469	0.3859
0.50	0.5898	0.6586	0.7984	0.8770	0.3362

7 DISCUSSION

Write your discussion here. Do not forget to use sub-sections. Normally, the discussion starts with comparing your results to other studies as precisely as possible. The limitations should be reflected upon in terms such as reproducibility, scalability, generalizability, reliability and validity. It is also important to mention ethical concerns.

8 CONCLUSION

Write your conclusion here. Be sure that the relation between the research gap and your contribution is clear. Be honest about how limitations in the study qualify the answer on the research question.

REFERENCES

- [1] Yasser Alginahi. 2010. *Preprocessing Techniques in Character Recognition*. <https://doi.org/10.5772/9776>
- [2] Sofia Ares Oliveira, Benoit Seguin, and Frederic Kaplan. 2018. dhSegment: A Generic Deep-Learning Approach for Document Segmentation. In *2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR)*. IEEE. <https://doi.org/10.1109/icfhr-2018.2018.00011>
- [3] Marius Bulacu, Rutger van Koert, Lambert Schomaker, and Tijn van der Zant. 2007. Layout Analysis of Handwritten Historical Documents for Searching the Archive of the Cabinet of the Dutch Queen. In *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, Vol. 1. 357–361. <https://doi.org/10.1109/ICDAR.2007.4378732>
- [4] Tieming Chen, Guangyuan Fu, Hongqiao Wang, and Yuan Li. 2020. Research on Influence of Image Preprocessing on Handwritten Number Recognition Accuracy. In *The 8th International Conference on Computer Engineering and Networks (CENet2018)*, Qi Liu, Mustafa Misir, Xin Wang, and Weiping Liu (Eds.). Springer International Publishing, Cham, 253–260.
- [5] Sergio Correia and Stephan Luck. 2023. Digitizing historical balance sheet data: A practitioner's guide. *Explorations in Economic History* 87 (2023), 101475. <https://doi.org/10.1016/j.eeh.2022.101475> Methodological Advances in the Extraction and Analysis of Historical Data.
- [6] Ahmad Drobj, Berat Kurar Barakat, Reem Alaasam, Boraq Madi, Irina Rabaev, and Jihad El-Sana. 2022. Text Line Extraction in Historical Documents Using Mask R-CNN. *Signals* 3, 3 (2022), 535–549. <https://doi.org/10.3390/signals3030032>
- [7] Wafaa S. El-Kassas, Cherif R. Salama, Ahmed A. Rafea, and Hoda K. Mohamed. 2021. Automatic text summarization: A comprehensive survey. *Expert Systems with Applications* 165 (2021), 113679. <https://doi.org/10.1016/j.eswa.2020.113679>
- [8] Sotirios Kastanas, Shaomu Tan, and Yi He. 2023. Document AI: A Comparative Study of Transformer-Based, Graph-Based Models, and Convolutional Neural Networks For Document Layout Analysis. arXiv:2308.15517 [cs.CL]
- [9] Constantin Lehenmeier, Manuel Burghardt, and Bernadette Mischka. 2020. *Layout Detection and Table Recognition – Recent Challenges in Digitizing Historical Documents and Handwritten Tabular Data*. Springer International Publishing, 229–242. https://doi.org/10.1007/978-3-030-54956-5_17
- [10] Minghao Li, Tengchao Lv, Jingye Chen, Lei Cui, Yijuan Lu, Dinei Florencio, Cha Zhang, Zhoujun Li, and Furu Wei. 2022. TrOCR: Transformer-based Optical Character Recognition with Pre-trained Models. arXiv:2109.10282 [cs.CL]
- [11] Xusheng Liang, Abbas Cheddad, and Johan Hall. 2021. Comparative Study of Layout Analysis of Tabulated Historical Documents. *Big Data Research* 24 (2021), 100195. <https://doi.org/10.1016/j.bdr.2021.100195>
- [12] Huaming Liu, Xuehui Bi, and Weilan Wang. 2019. Layout analysis of historical Tibetan documents. In *2019 2nd International Conference on Artificial Intelligence and Big Data (ICAIBD)*. 74–78. <https://doi.org/10.1109/ICAIBD.2019.8837040>
- [13] Tengchao Lv, Yupan Huang, Jingye Chen, Lei Cui, Shuming Ma, Yaoyao Chang, Shaohan Huang, Wenhui Wang, Li Dong, Weiyao Luo, Shaoxiang Wu, Guoxin Wang, Cha Zhang, and Furu Wei. 2023. Kosmos-2.5: A Multimodal Literate Model. arXiv:2309.11419 [cs.CL]
- [14] Siddharth Misra and Yaokun Wu. 2020. Chapter 10 - Machine learning assisted segmentation of scanning electron microscopy images of organic-rich shales with feature extraction and feature ranking. In *Machine Learning for Subsurface Characterization*, Siddharth Misra, Hao Li, and Jiaobo He (Eds.). Gulf Professional Publishing, 289–314. <https://doi.org/10.1016/B978-0-12-817736-5.00010-7>
- [15] Arthur Flor de Sousa Neto, Byron Leite Dantas Bezerra, and Alejandro Héctor Toselli. 2020. Towards the Natural Language Processing as Spelling Correction for Offline Handwritten Text Recognition Systems. *Applied Sciences* 10, 21 (2020). <https://doi.org/10.3390/app10217711>
- [16] OpenCV. 2024. Morphological Transformations. https://docs.opencv.org/3.4/dd/d7/tutorial_morph_lines_detection.html Accessed: 2024-04-18.
- [17] Rafael Padilla, Sergio Netto, and Eduardo da Silva. 2020. A Survey on Performance Metrics for Object-Detection Algorithms. <https://doi.org/10.1109/IWSSIP48289.2020>
- [18] Rémi Petitpierre, Marion Kramer, and Lucas Rappo. 2023. An end-to-end pipeline for historical censuses processing. *International Journal on Document Analysis and Recognition (IJDAR)* 26, 4 (March 2023), 419–432. <https://doi.org/10.1007/s10032-023-00428-9>
- [19] Jose Ramón Prieto, José Andrés, Emilio Granell, Joan Andreu Sánchez, and Enrique Vidal. 2023. Information extraction in handwritten historical logbooks. *Pattern Recognition Letters* 172 (2023), 128–136. <https://doi.org/10.1016/j.patrec.2023.06.008>
- [20] Stephen Quinn and William Roberds. 2009. *An economic explanation of the early Bank of Amsterdam, debasement, bills of exchange and the emergence of the first central bank*. Cambridge University Press, 32–70.
- [21] Mathias Seuret. [n.d.]. *Layout Analysis in Handwritten Historical Documents*. Chapter Chapter 4, 45–65. https://doi.org/10.1142/9789811203244_0004 arXiv:https://www.worldscientific.com/doi/pdf/10.1142/9789811203244_0004
- [22] Brandon Smock, Rohith Pesala, and Robin Abraham. 2022. PubTables-1M: Towards Comprehensive Table Extraction From Unstructured Documents. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 4634–4642.
- [23] Enrique Vidal, Alejandro H. Toselli, Antonio Ríos-Vila, and Jorge Calvo-Zaragoza. 2023. End-to-End page-Level assessment of handwritten text recognition. *Pattern Recognition* 142 (2023), 109695. <https://doi.org/10.1016/j.patcog.2023.109695>
- [24] Yiheng Xu, Tengchao Lv, Lei Cui, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, and Furu Wei. 2021. LayoutXML: Multimodal Pre-training for Multilingual Visually-rich Document Understanding. arXiv:2104.08836 [cs.CL]

525 Appendix A SUPPLEMENTARY FIGURES

61		62	
<i>Gan Van Os d. 8.</i>		<i>Ced. ady 4 Aug 1684</i>	
Buy. van Crim Potte	1227 -	Van. Mondt	15768 17
Mr. Bontje	963 2000	van Blizting	355 2000
Mr. Dordille	226 592 10	Gillerty & Sandy N.Y. Gathering	326 2800
Mr. E. Parus	120 497 9	Mr. & Mrs. T. & S. Lingard	168 289 14 -
Mr. G. Verloet @ Wrom	33 1947 19	Bishop's Lodge	820 462 -
Mr. G. Verloet @ Wrom	1184 9710	Mr. Athur	988 2000
Mr. G. Verloet @ Wrom	5285	Mr. et Mrs. C. K. H.	867 1163 19 0
Brueghel Graeber	1126 850	Dick Van Domburg	240 3000
Graeber Graeber	1016 1500	Dirk Van Domburg	862 1400
Graeber Graeber	959 1635	Mar. & Mr. Van Brumhoff	301 1000
Graeber Graeber	900 1075	Baner	30958 10 8
Graeber Graeber	313 1630 15	Dars. Lep. 6	179 1200
Mr. Bontje	225 1140	Poel, B. 3000	122 3000
Mr. Bontje	139 1140	Pan Yab. V. 1000	131 1320 -
Car. & Sons Gerard	15 2729 16 0	Mr. Nees Maer	840 1000
Dick Van Domburg	302 600	Louis Garinier	202 2029 18 8
Van. Mayne	314 1200	Van Rosloff & Van. Werft	39036 9 -
Jos. Van Gertach	19318 11 0	God. Meylandt	311 1000
Jos. Van Gertach	1099 2518 15	Sam. Collier	122 1500
Van. f. d. wafay	877 3108 18	Mr. C. Car. & S. Gerard	782 6000
Van. Mondt	315 809 17	Van. Mondt	11 950
Mr. Bontje	225 1200	Van. 35125 1000	355 2000
Van. Bontje	121 1340 12 0	Van. 35125 1000	800 3000
Van. Weller	1148 28288 9 0	Van. 35125 1000	81 123 15
Van. Weller	1148 350 9 0	Van. 35125 1000	122 2000
Van. Weller	1128 1000	Van. 35125 1000	908 905 -
Van. Weller	1039 1000	Held. D. Van. G. G. West	820 6425 4 -
Van. Weller	1012 2600	Held. D. Van. G. G. West	1011 1000
Van. Weller	869 9 9	Held. D. Van. G. G. West	1011 1525 10 -
Samuel Eliza Schra Land	33507 13 8	Baude. Charles	355 300
Matz Etteman	113 960	Car. & K. H.	816 1500
John Gorius	107 1200	if. Minifart & Elckhoff	973 504 2 0
Jan. Son. Minifart	121 900	Van. Van. Zelln	291 15000
Mr. Arout	103 2558		84513 18
Mr. Arout	900		
Mr. Arout	90025 13 0		
Mr. Arout	802 300		
Flayo Garlick	979 800		
Mr. Arout	44 1600 0 -		
Mr. Arout	43 2350		
Mr. Arout	1087 600		
Mr. Arout			
12. Lambert ten Cate. Ob. Kade	45681 1 8		
13. Van. Collaer	116 1000		
14. Van. Mayne	179 900		
15. Van. Fontain Jan	314 1100		
16. C. V. C. J. H. H.	108 1300		
17. G. G. G. H.	814 126		
18. G. G. G. H.	363 600		
19. S. J. H. H.	252 1854 1		
Jan. S. H. H.	268 1060		
Dick Van Domburg	302 1000		
Das. Balde	- 985 2400		
Das. Garlick			
Das. Garlick	979 7100 2 8		
Das. Garlick	992 1000		
Das. Garlick	992 1000		
Das. Garlick	425 1500		
Das. Garlick	372 1600		
Das. Garlick	939 839 0		
Das. Garlick	63660 10 8		
Das. Garlick	895 1400		
Das. Garlick	310 1011 13 0		
Das. Garlick	252 783 1		
Das. Garlick	119 1000		
Das. Garlick	62 756 5		
Das. Garlick	437 16575 2 8		
Das. Garlick	84513 1 8		

Figure 9: An example page from the ledger collection.

(a) The PageXML data of our ground truth plotted on part of the original image.

(b) The PageXML data of the HTR system plotted on part of the original image.

Figure 10: The PageXML data of the ground truth and HTR system plotted on the original image. The tabular structure of the ground truth values can be seen where this is not the case in the HTR output.

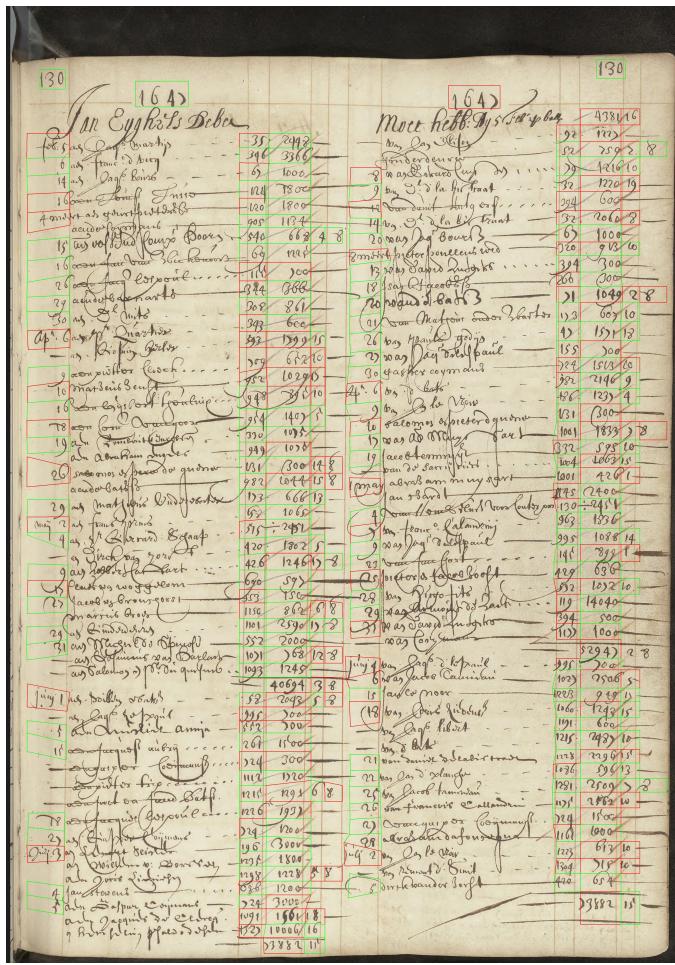


Figure 11: An example page showing the words correctly matched by the HTR system in green and the words wrongly matched in red. This can be helpful in visually inspecting the page to see whether there are common characteristics on parts of the page that could cause mistakes.

526 Appendix B DETAILED METHODOLOGY

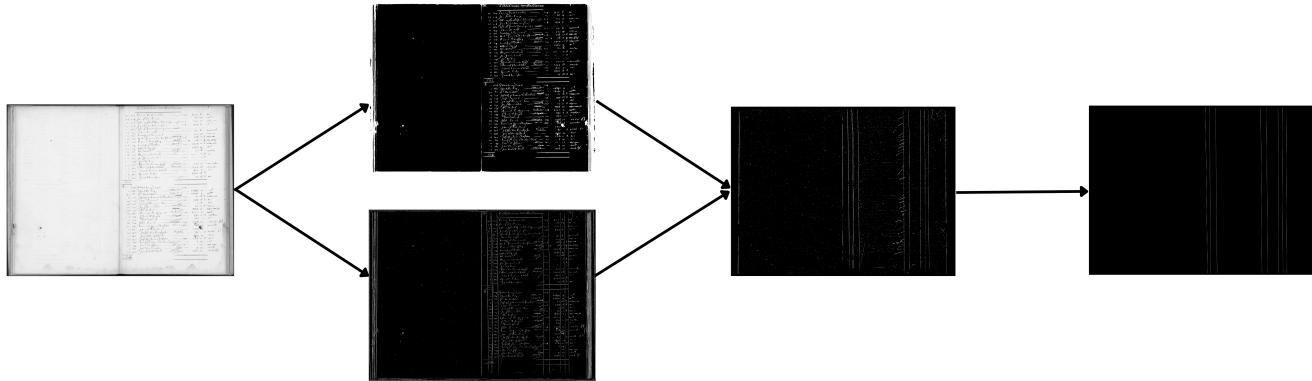


Figure 12: The process of the extraction of column lines. We start with the blurred grayscale image. From there we create two threshold images, one containing the faint column lines and one that does not contain them. We subtract the images from each other, remove noise and use Hough line detector to find the columns lines.



Figure 13: The process of extracting the rows from the ledger page. We use the HTR system to determine bounding boxes of words from the original page. We then determine the row bounding boxes based on the y-coordinates of the word bounding boxes.