

INSTITUT FÜR INFORMATIK

A Systematic Comparison of Different Approaches for Unsupervised Extraction of Text from Scholarly Figures [Extended Report]

Falk Bösch
Ansgar Scherp
Bericht Nr. 1607
November 2016
ISSN 2192-6247



CHRISTIAN-ALBRECHTS-UNIVERSITÄT
ZU KIEL

Institut für Informatik der
Christian-Albrechts-Universität zu Kiel
Olshausenstr. 40
D – 24098 Kiel

**A Systematic Comparison of Different
Approaches for Unsupervised Extraction of Text
from Scholarly Figures [Extended Report]**

Falk Bösch
Ansgar Scherp

Bericht Nr. 1607
November 2016
ISSN 2192-6247

e-mail: fboe@informatik.uni-kiel.de,
asc@informatik.uni-kiel.de

An abridged version of this work is published at the 23rd International
Conference on Multimedia Modeling, Reykjavik, Iceland, January 2017.

A Systematic Comparison of Different Approaches for Unsupervised Extraction of Text from Scholarly Figures

[Extended Report]

Falk Böschchen¹ and Ansgar Scherp^{1,2}

¹ Kiel University, Kiel, Germany
{fboe,asc}@informatik.uni-kiel.de

² ZBW - Leibniz Information Centre for Economics, Kiel, Germany
a.scherp@zbw.eu

Abstract. Different approaches have been proposed in the past to address the challenge of extracting text from scholarly figures. However, so far a comparative evaluation of the different approaches has not been conducted. Based on an extensive study, we compare the 7 most relevant approaches described in the literature as well as 25 systematic combinations of methods for extracting text from scholarly figures. To this end, we define a generic pipeline, consisting of six individual steps. We map the existing approaches to this pipeline and re-implement their methods for each pipeline step. The method-wise re-implementation allows to freely combine the different possible methods for each pipeline step. Overall, we have evaluated 32 different pipeline configurations and systematically compared the different methods and approaches. We evaluate the pipeline configurations over four datasets of scholarly figures of different origin and characteristics. The quality of the extraction results is assessed using F-measure and Levenshtein distance. In addition, we measure the runtime performance. The experimental results show that there is an approach that overall shows the best text extraction quality on all datasets. Regarding runtime, we observe huge differences from very fast approaches to those running for several weeks.

Keywords: Scholarly Figures · Text Extraction · Comparison

1 Introduction

Scholarly figures are data visualizations in scientific papers such as bar charts, line charts, and scatter plots [9]. Different research has been conducted to extract and use text from figures like translating the text to Braille [17], re-engineering the raw data from the figures [28], or just for image search [31]. Many approaches follow a semi-supervised text extraction approach [8, 28]. However, semi-supervised approaches do not scale with the amount of scientific literature published today. Thus, unsupervised methods are needed to address the task of text extraction from scholarly figures. This task is challenging due to the heterogeneity in the appearances of the scholarly figures such as varying colors, font

sizes, and text orientations. Nevertheless, extracting text from scholarly figures, such as the examples shown in Figure 1, is an important task as the text provides additional information that is not contained in the papers [3]. To the best of our knowledge, no comparison of the different approaches for text extraction from scholarly figures has been conducted so far. Most likely the reason behind this is that existing works come from various different research areas. Furthermore, there are no annotated, high-quality datasets as gold standard publicly available. The only exception is the CHIME dataset. It has high quality gold standard annotations, but it does not contain orientation information. In addition, several of the scanned figures are only provided in low-quality resolution. Thus, a thorough assessment and comparison of the different text extraction approaches has not been conducted.

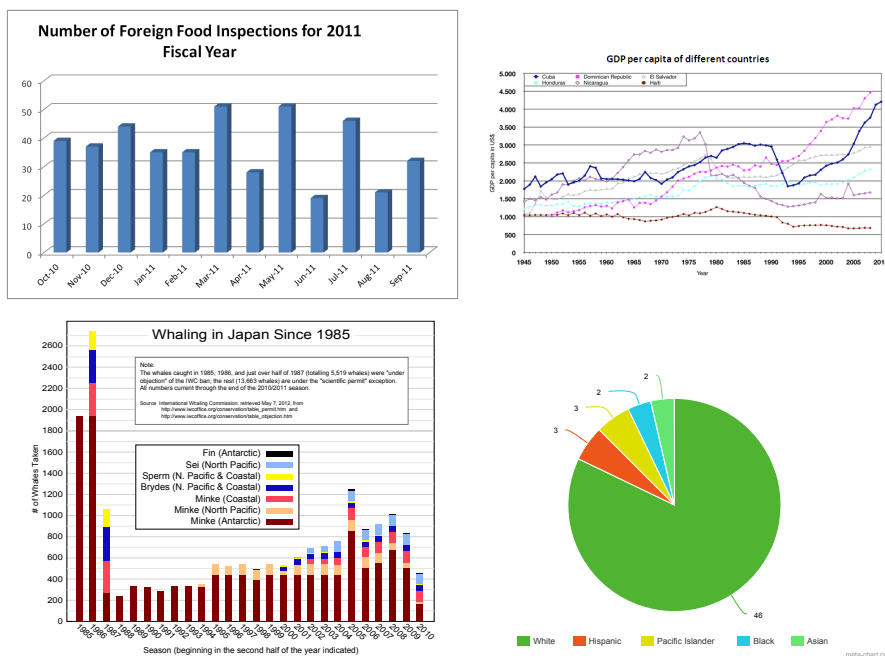


Fig. 1. Exemplary scholarly figures (Source: Wikimedia Commons (Public Domain))

Based on the related work, we have defined a generic pipeline of six sequential steps that abstracts from the various works on text extraction from scholarly figures. We have re-implemented and systematically evaluated the most relevant approaches for text extraction from scholarly figures as described in the literature. In total, we have investigated 7 different pipeline configurations motivated by approaches described in the literature. Each configuration is a combination of six to nine methods for the sequential steps in the extraction pipeline. Fur-

thermore, we have created 25 modifications of the best performing pipeline configuration to systematically measure the influence of different methods applied in the text extraction pipeline. Thus, in summary we have compared 32 different configurations for text extraction from scholarly figures in this paper.

We assess each pipeline configuration with regard to the accuracy of the text location detection via precision, recall, and F1-measure. In addition, we evaluate the text recognition quality using Levenshtein distance. Finally, we are comparing the runtime of the different configurations to find the most efficient ones. We use four datasets in our evaluation: one from economics [2] (EconBiz), one synthetically generated [18] (CHIME-S), one scanned and collected on the Internet [32] (CHIME-R), and one created from figures of academic books provided by the publisher DeGruyter³ (DeGruyter). We manually labeled the EconBiz dataset and DeGruyter dataset, while the CHIME datasets were created in 2006 by the Center for Information Mining and Extraction, School of Computing, National University of Singapore.

In summary, the contributions of the paper are as follows:

- (i) We conduct a systematic comparison of in total 32 configurations of a generic pipeline for text extraction from scholarly figures. Each configuration consists of a combination of six to nine methods from a total of 21 different methods that we have implemented and evaluated.
- (ii) We make available four manually labeled datasets of scholarly figures that allow reproducing and extending our results.⁴
- (iii) Furthermore, we also make available the implementation of our generic pipeline, including the 21 methods, as well as the 32 configurations for text extraction from scholarly figures that we compared. This allows for reproducing our results and applying our generic implementation on other datasets.

The subsequent section discusses the related work in the field. It serves as foundation for our generic pipeline and the different methods used in the pipeline configurations. In Section 3, the sequence of six processing steps of the generic pipeline is briefly described. The individual methods that are selected for comparison in this work are described in Section 4 and the pipeline configurations are introduced in Section 5. Section 6 describes the four datasets and the measures used in our evaluation. The results are described in Section 7 and discussed in Section 8, before we conclude.

2 Related Work

Text extraction from scholarly figures is addressed by research groups from different domains. Thus, one finds different terms that basically describe the same concept, such as information graphics [3], figures [9], charts [15], diagrams [5],

³ <http://www.degruyter.com/>

⁴ <http://www.kd.informatik.uni-kiel.de/en/research/software/text-extraction>

and different variations of them. In the following, we will commonly denote them as scholarly figures or short figures.

First, we discuss relevant works on text extraction from scholarly figures. Subsequently, we consider cartographic maps, domain-specific approaches from life sciences and chemistry, as well as briefly discuss approaches for making figures accessible to visually impaired users as well as applying text extraction on natural photos.

Scholarly Figures An early work on text extraction from scholarly figures is by Huang et al. [14, 15]. Their text extraction pipeline starts with a Connected Component Labeling (CCL) [26] that generates components of coherent text elements and graphics elements. In a subsequent step, these elements are separated by applying a series of filters. In the next step, the text elements are grouped using a derivation of Newton’s formula for “Gravity” from classical physics. The authors claim that this method is capable of separating text at different orientations into different groups of text elements. Optical character recognition (OCR) is applied on these text groups and the recognized text is classified into strings and numbers. Finally, the recognized text is manually corrected in order to have a clean assignment to the corresponding graphical elements. Sas and Zolnierak [27] propose a three-stage approach for text extraction from figures. Starting with a conversion of the input image to gray-scale, the authors apply filtering operations, binarization, and CCL to generate coherent regions. Regions are filtered by empirical thresholds and classified into text and graphic elements using a decision tree. Tesseract is used for OCR. Besides normal orientation, the input to the OCR engine is also rotated at a 90° angle to capture vertical text elements such as labels of the y-axis. Finally, the text detection is verified by assessing the number of special characters recognized in the text regions. Unfortunately, the authors did not assess the quality of their OCR results. DiagramFlyer [5] is a retrieval system for bar charts or scatter plots. Although the approach sounds to be very promising, the authors do not provide any technical specifications or details on how their extraction pipeline works. This renders it impossible to consider their work in a systematic comparison. Finally, we have developed a pipeline called TX for unsupervised text extraction from scholarly figures [1, 2]. The TX pipeline uses an adaptive binarization method based on Otsu’s method. Subsequently, CCL is applied to extract coherent regions. A few heuristic rules are applied before the regions are clustered using DBSCAN in order to separate text elements from graphical elements. A Minimum Spanning

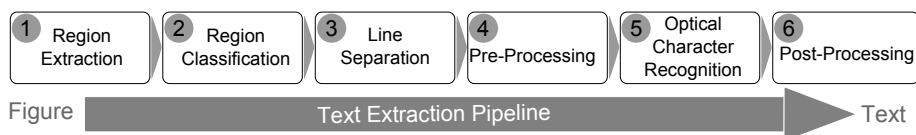


Fig. 2. Generic pipeline for text extraction from figures abstracted from the literature

Tree (MST) clustering is applied to detect single text lines. The orientation of these text lines is computed using a discrete Hough transformation and each line is rotated into horizontal mode in order to send it to a standard OCR engine. Here, the Tesseract⁵ OCR engine is used.

Cartographic Maps Cartographic maps use text elements to show city and street names, regions, and landmarks. An early work on text extraction from maps is the approach by Deseilligny et al. [11] which relies on CCL for extracting regions. However, in contrast to the works on scholarly figures discussed above, Deseilligny et al. normalize each region in order to apply a rotation invariant character recognition. Multiple character hypotheses are generated for each region and those hypotheses are selected which create coherent strings and follow specific syntactic rules.

A more recent approach is the semi-automatic text extraction proposed by Chiang et al. [8]. In contrast to most of the other works, the input image is not converted to gray-scale. Instead, a color quantization algorithm is applied. The authors separate text elements from graphical elements using a run-length smoothing algorithm based semi-automatic extraction that requires a positive and a negative example for each text-color/background-color combination. Text lines are detected by applying dilation operators on the connected components. The orientation of each line is estimated using a Single String Orientation Detection algorithm, which is based on morphological operations. The algorithm evaluates all possible orientations of text line candidates in a brute-force manner. The text line is rotated to horizontal orientation and AbbyyFineReader⁶ is applied for OCR. After the OCR phase, a recognition confidence score is computed to filter the results.

Domain Specific Text Extraction A text detection algorithm for biomedical images was proposed by Xu and Krauthammer [31] as part of the Yale Image Finder. The authors first detect and remove so-called layout elements, followed by a binarization, median filter, and edge detection with the Sobel operator. The text region extraction, based on horizontal and vertical histogram projection analysis, is conducted on the edge image. This is performed recursively until the image cannot be split any further. During this recursive processing of the regions, heuristic filters are applied to only subdivide those regions that contain text and discard the others.

Lu et al. [23] developed a retrieval engine for scholarly figures in chemistry. First, the input image is converted to gray-scale and an edge image is computed. A Hough transformation is applied on the edge image in order to compute a feature vector. The feature vector is used to classify the input in order to find 2D plots. Only 2D plots are further processed. First, the 2D plot is binarized. Then, the axes are detected and the plot is segmented by applying CCL. The text de-

⁵ <https://github.com/tesseract-ocr/>

⁶ <http://www.abbyy.com/ocr-sdk/>

tection is based on fuzzy rules and includes a method for separating overlapping characters. The recognition of text strings is conducted using GOCR⁷.

Access for the Visually Impaired Another approach that requires text extraction from figures is the work by Jayant et al. [17]. Their goal is to translate figures into Braille language for the visually impaired. First, a color reduction is conducted with Adobe Photoshop. Subsequently, the figure is manually classified into a set of predefined figure types. CCL is applied to the figure to extract regions. In order to separate text elements from graphical elements, the authors manually train a Support Vector Machine (SVM) per figure type as well as per book where the figures were taken from. Thus, the authors make the assumptions that all figures of a certain type have a similar design throughout a single book. Subsequently, a separation into text line structures is performed, using a so-called label training algorithm which uses a Minimum Spanning Tree with manually created test data. The text line orientation is estimated by minimizing the perpendicular squared distance. Finally, OCR is conducted with Omnipage⁸ or AbbyyFineReader. Carberry et al. [4] analyze figures, especially bar charts, pie charts, and line charts, to generate textual summaries for visually impaired users. Their Visual Extraction Module (VEM) claims to be capable of extracting text elements and their position [6]. However, the paper does not provide technical details on how this is achieved. In their current work [4], manually generated datasets are used instead of the VEM, which may indicate that the VEM does not generate output of sufficiently high quality.

Approaches on Natural Photos Besides text extraction from scholarly figures and related images, there is also research regarding OCR on natural photos. The research in this area has several interesting ideas to solve the OCR problem. For example Olszewska [24] presented a template-matching approach to extract numbers of arbitrary position and orientation in the captured 3D space in images using contour information. Other approaches show promising results as well [12, 22]. However, the works for text extraction from natural photos often make specific assumptions about the difference in appearance of text and background/graphic elements. For example, the assumptions that text elements are generally smaller than graphic elements [13], that text elements can be identified via their edges [22], or that text has a unique or more homogeneous color [13]. These assumption often do not hold for scholarly figures like charts, diagrams, or graphics.

3 Generic Pipeline

Based on the related work, we derived a generic pipeline for text extraction from scholarly figures. The pipeline consists of six sequential steps as illustrated in Figure 2. Each step can be implemented by different methods. This allows to

⁷ <http://www-e.uni-magdeburg.de/jschulen/ocr/index.html>

⁸ <http://www.nuance.com/for-business/by-product/omnipage/csdk/index.htm>

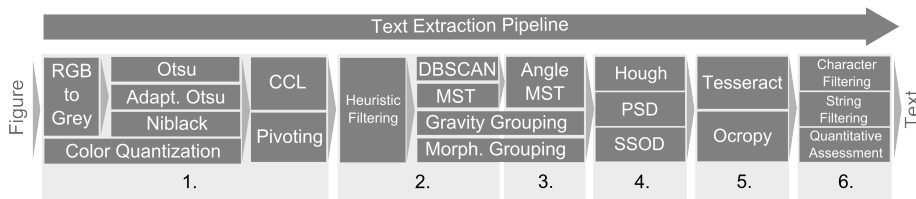


Fig. 3. Overview of the possible pipeline configurations

compare the pipeline configurations from the literature as well as to provide new configurations of so far not investigated combinations of the various methods. Below, we provide a brief summary of the different sequential steps that we identified for the generic pipeline. A formalization of the pipeline can be found in our earlier publication [1]. Please note, for describing the steps of our generic pipeline we use the following terminology: We refer to scholarly figures as **images**, since it is the accepted term in computer vision. A **region** is a set of pixels of an image. Each region constitutes either one or sometimes multiple text characters or graphical symbols. A **text line** or text element is a set of regions representing text.

The input to the pipeline is a (color) image of a scholarly figure and the output are text elements together with their position, dimension, and orientation. The six steps are as follows: **(1)** The first step extracts regions from an image. Thus, a decision on pixel-level has to be made about what part of the image belongs to a region and what is background. A common algorithm for this task is Connected Component Labeling (CCL) on a binarized image. **(2)** In the second step, the previously computed regions are classified either as text or graphics. The regions classified as graphics are ignored in the subsequent steps. **(3)** The third step computes text lines from the text regions provided by the previous step. It is necessary to compute the text lines since most OCR engines work only on horizontal text input and one can only reliably estimate the orientation of single lines of text. **(4)** The fourth step estimates the orientation of the text lines and performs other pre-processing that a specific OCR engine might need. Besides rotating text lines to horizontal orientation, one may need to scale text lines to a sufficiently high resolution or remove noise. **(5)** Subsequently, the fifth step actually performs the Optical Character Recognition (OCR). A commonly used OCR engine is Tesseract, developed by Google and used in the Google Books project. **(6)** Finally, post-processing is applied on the OCR results. For example, the OCR output is corrected using some heuristics.

4 Methods for the Pipeline Steps

For each step of the generic pipeline, we compare different methods motivated by approaches described in the literature. Below, we describe the methods selected

along the steps of the generic pipeline as shown in Figure 2. The large number of methods allow only a brief description of each method, but further details can be found in the references. In Section 5, we discuss different configurations of the generic pipeline assembled from the methods described below. An overview of the possible configurations can be found in Figure 3.

4.1 Step (1): Region Extraction

The region extraction step consists of two sub-steps: First, the input is transformed from color space into one or multiple binary images as described below. Second, the actual region extraction is conducted using one of two approaches that we identified in the literature. Overall, the output of this step is a set of binary regions where each region represents one or more text characters or graphical symbols.

Binarization of Color Images Given a color image, one can convert it directly to multiple binary images. An alternative is to use an intermediate transformation to a grey-scale image, which is then converted to the binary output. For directly converting a color image to multiple binary images, a so-called **Color Quantization** method can be used. Color Quantization performs a clustering over the color space. For each cluster, a representative color is chosen and each color in the original image is replaced by the color of the cluster closest to it. Finally, the image is split into multiple binary images, one for each color, where all pixels that have the specific color are set and the others are not. The color quantization method is inspired by the work of Chiang [7], Fraz [12], and Jayant [17]. The **RGB to Grey via Luminance** method is commonly used to create an intermediate grey-scale image [2, 14]. It uses the conversion formula $Y = 0.2126R + 0.7152G + 0.0722B$ to weight the color components red (R), green (G), and blue (B) to the luminance value Y based on human perception. Subsequently, the grey-scale image is binarized using one of the following methods: **Otsu’s Method** [25] separates two regions by finding the threshold that maximizes the inter-class variance. One problem with Otsu’s method is that it only computes one threshold (or multiple if the Multi-Otsu method is used) which are applied globally on the entire image. This leads to problems with local inhomogeneities like varying text-color/background-color combinations. In order to address this challenge, an **Adaptive Otsu Binarization** method [2] was developed that computes multiple thresholds to create a locally adaptive binarization. This is achieved by subdividing the image into smaller parts and recursively applying Otsu’s method to compute new thresholds. Another binarization method is **Niblack’s Method** and modifications of it [21]. We tested all versions of Niblack’s Method as described by Khurshid et al. [21] with the specified parameters. The modified version of Sauvola performed best during our preliminary tests. Thus, we compare it with the two Otsu variants.

Region Extraction from Binary Images Given a binary image, we can apply one of the following two methods for region extraction: The most common method [2,

14, 27] is **Connected Component Labeling (CCL)** [26]. The CCL algorithm iterates over the whole image and assigns each pixel to a region by taking the assignment of the adjacent pixels into account (4-/ 8-pixel-neighborhood). As an alternative, we use Xu and Krauthammer’s [31] **Pivoting Histogram Projection** method for region extraction. Here, first an edge image of the binary image is computed. The edge image’s pixels are alternately projected on the x- and y-axes and the image is split after every projection at the minimal point(s) of the histogram into multiple sub-images. Each sub-image is further processed while alternating the direction until no further split is possible. Thus, one obtains multiple rectangular areas, which are converted into regions by taking all foreground pixels of each region from the binary image.

4.2 Step (2): Region Classification

The output of the region extraction in Section 4.1 has to be classified into text elements and graphical elements. For each region, we compute a feature vector which is composed of the center of mass x-/y-coordinates, width and height, and area-occupation-ratio (the number of foreground pixels of the bounding box divided by the area) that are used to group the regions into text and graphics. **Heuristic Filtering** methods are commonly applied as pre-processing to help separate text elements from graphical elements [2, 27]. For example, very small regions typically constitute noise, large regions are graphical elements like axes, and average-sized regions refer to characters. Another approach is to consider the coverage of the bounding box of a region. Heuristic filters are parameterized and thus require a suitable choice of parameter values. A less parameterized version for region classification are unsupervised density-based clustering algorithms like **DBSCAN** which can be used to group regions [2]. Since text is normally more dense and of different size than graphic elements, it can be separated from graphic elements using DBSCAN. Another method to distinguish between text and graphics is the graph-based **Minimum Spanning Tree (MST)** clustering algorithm [17]. After constructing the tree, the clusters are created by splitting the graph, i. e., removing edges. There are multiple possibilities on how to split the tree. In our experiments, we consider splits at inconsistent edges, i. e., edges that are longer than the local average edge length. Huang et al. [15] proposed **Grouping Rules based on Newtons Gravity Formula** from classical physics. The formula computes a threshold which defines whether two regions should be grouped together. This results in text elements, each representing a single text line. Finally, the **Morphological Method** by Chiang et al. [3, 8] is applied on the individual pixels of the image and uses morphological operations to merge characters into words. Thus, it also generates text elements.

4.3 Step (3): Separation into Text Lines

The methods for determining text elements in the previous step can result in clusters of regions of various shape. Thus, the text elements may contain text of different orientation. This is problematic for standard OCR engines as they

require text at horizontal orientation. Therefore, it is necessary to split the text elements into single text lines, since one can only reliably estimate the orientation for single lines of text. One method to separate text elements into text lines is to apply an **Angle-Based MST** [2] clustering. The MST is constructed over the centers of mass of the regions in a cluster. The assumption is that characters of a text line are closer to each other than characters from different text lines. Thus, most edges of the MST constitute a single line while only few edges connect across different text lines. Edges connecting different lines will have orientations that differ strongly from the main orientation and can therefore easily be removed.

4.4 Step (4): Text Line Orientation

Standard OCR engines require that the text of the input image has horizontal orientation. In this step, the single text lines produced in the previous step are analyzed w.r.t. their orientation. We compute the orientation of each line using one of the following methods: The first method uses the **Hough Transformation** [16] to calculate the orientation [2]. The method transforms all center of mass coordinates of the characters of a text line into Hough space and the maximal value in this Hough space represents the main orientation. Since text lines can have only an orientation between +90 and -90 degree, we can limit the Hough space computation to this interval. Another method to estimate the orientation is to minimize the **Perpendicular Squared Distance** of the bounding box of each text line [17]. The **Single String Orientation Detection (SSOD)** [8] method assesses different orientation candidates by rotating the text elements and applying morphological operations on its regions. The orientation at which the largest pixel area remains after applying the operators, is selected as the orientation of the text element. Subsequently, we rotate the text lines into the opposite direction to bring them into horizontal alignment.

4.5 Step (5): Optical Character Recognition

We now have individual text lines at horizontal orientation. Thus, in this step, we can apply standard Optical Character Recognition (OCR) engines to extract the text. We analyzed different OCR engines mentioned in the related work. We have selected Tesseract and Ocropy as they are freely available and frequently updated. A possible extension is the use of commercial OCR engines like AbbyFineReader. To ensure the reproducibility of the results of this work, we decided to limit the number of options for the OCR step to open source solutions. The OCR engine **Tesseract** provides trained models for different languages, where we choose English. We are not using Tesseracts Layout Analysis capabilities for removing graphic elements and conducting line detection since it only works for horizontal text. Tesseract has already been used in the past for text extraction from scholarly figures [2, 27]. **Ocropy**⁹ is a collection of open source tools for document analysis and OCR. It is designed for character recognition

⁹ <https://github.com/tmbdev/ocropy>

from full-page documents similar to Tesseract. Ocropy has several constraints regarding parameters like minimum image width and height in order to assure good results. Thus, we modify the input image to fulfill the required thresholds by properly scaling the text lines. Like Tesseract, the OCR engine Ocropy provides a pre-trained model for the English language.

4.6 Step (6): Post-Processing

The last step of the pipeline conducts potential corrections of the textual output of OCR engines. Here, several approaches exist: For example, one can use dictionaries [29] to correct the OCR output. Since text in scholarly figures is very sparse and often contains abbreviations, one cannot apply standard dictionaries. A much simpler heuristic-based correction is the **Special Character Filtering per single Character** method. It removes all special characters from the output, i. e., all characters that are not a white space, number, or character from a-z or A-Z. This makes sense, because recognition errors often appear in form of special characters like dots or dashes in the output. Sas and Zolnierek proposed a **Special Character Filtering per String** [27], which is a modified version of the previous method. Here, complete text elements are removed if they contain too many special characters. Another post-processing method is the **Quantitative OCR Assessment** by Jayant et al. [17]. The main idea is to reuse knowledge from previous steps of the extraction pipeline by comparing the number of regions that went into the OCR process with the number of characters that were recognized from them. If the difference is above a certain threshold, one can assume that one or multiple recognition errors happened during the OCR process. While Jayant’s approach also takes recognition confidence information on character level from AbbyyFineReader into account, one cannot in general assume that this information is available from all OCR engines. Thus, the implemented **Quantitative OCR Assessment** method performs its post-processing only based on the difference between the number of regions before and the number of characters after the OCR step.

5 Pipeline Configurations

From the methods defined in the previous section, one can create various pipeline configurations. Some methods are restricted in how they can be combined as illustrated in Figure 3. Section 5.1 discusses the configurations that are motivated from the related work. A further systematic assessment of the different methods is conducted by the configurations described in Section 5.2.

5.1 Motivated from the Literature

There are seven pipeline configurations that are motivated from the literature. Each configuration is identified by **(x)**, an acronym created from the contributing author(s). The first configuration (**SZ13**) is inspired by the work of Sas

and Zolnierek [27]. It uses Otsu’s method for binarization, followed by CCL. Subsequently, it applies heuristic filtering similar to the original approach. The decision tree used by Sas and Zolnierek is replaced by the line generation approach based on MST. The rationale behind this is that a decision tree is a supervised method while MST is unsupervised. This configuration does not use any method for orientation estimation from step 4 of the pipeline, since the original work by Sas and Zolnierek does not have such a feature. Tesseract is used as OCR engine, since it was also used in the original paper. In the post processing step, all strings are removed that contain too many special characters. The second configuration (**Hu05**) is based on the work of Huang et al. [15]. After region extraction using Otsu binarization and CCL, the Heuristic Filter method is applied, and the regions are grouped using the Gravity method. Finally, the grouped regions are processed with Tesseract. Based on the work of Jayant et al. [17], configuration (**Ja07**) starts with Otsu’s method and CCL. Subsequently, it clusters the regions using a MST and approximates the orientation by minimizing the perpendicular squared distance. Text recognition is achieved by applying Tesseract. Different from the previous configurations, the fourth configuration (**CK15**) – inspired by Chiang et al. [8] – uses Color Quantization to generate multiple binary images, followed by a CCL. Subsequently, it applies heuristic filtering and Morphological Clustering on the regions. This step differs from the original paper, where the relevant color levels were manually selected. Thus, we assess all extracted binary images. The orientation of each cluster is estimated using the SSOD method, followed by Tesseract OCR, and quantitative post-processing. Similar to the previous pipeline configuration, the fifth configuration (**Fr15**), inspired by Fraz et al. [12] from the photo processing domain, starts with Color Quantization and CCL. The original approach uses a supervised SVM to form words, which we replaced with unsupervised methods from our methods set. The extracted regions are filtered and DBSCAN is applied, followed by a MST clustering into text lines. The orientation of the text lines is calculated using Hough method and the text is recognized using Tesseract. All configurations so far use CCL to extract regions. The sixth configuration (**XK10**), motivated by Xu and Krauthammer [31], uses the pivoting algorithm after binarization with adaptive Otsu. The regions are filtered using heuristics and grouped into lines using DBSCAN and MST. This differs from the original work, which only applied heuristic filtering to remove the graphic regions. The reason behind this is that the authors only aimed at finding text regions and not to recognize the text. Thus, we filled the rest of the pipeline steps with suitable methods. The orientation of each line is estimated via Hough and OCR is conducted with Tesseract. Finally, configuration (**BS15**) resembles our own work [2]. It uses adaptive Otsu for binarization and CCL for region extraction. Heuristic Filtering is applied on the regions and DBSCAN groups them into text elements. Text lines are generated using the angle-based MST approach and the orientation of each line is estimated via Hough transformation, before applying Tesseract’s OCR.

5.2 Systematic Modifications

In order to evaluate the influence of the individual methods, we chose the pipeline configuration **(BS15)** as basis for systematical modification, since it is the most recent development for the task of automatically extracting text from scholarly figures and showed the best performance of the seven configurations from the literature. The systematic modifications are organized along the six steps of the generic pipeline in Figure 2. Each of the systematic configurations has an identifier **(BS-XYZ)** based on the original configuration, where X is a number that refers to the associated pipeline step and YZ uniquely identifies the method. The systematically modified configurations are described below:

Modifications of Step (1): The binarization and region extraction is evaluated with the following configurations: **(BS-1NC)** differs from **(BS15)** by using Niblack instead of adaptive Otsu for binarization. Configuration **(BS-1OC)** uses the third option for binarization, Otsu’s method. Color quantization is combined with the pivoting region extraction in **(BS-1QP)**.

Modification over Steps (2) and (3): The next step is the region classification and generation of text lines. Configuration **(BS-2nF)** differs from the base configuration by not applying the optional heuristic filtering method. Configuration **(BS-2CG)** uses the Gravity Grouping instead of DBSCAN and MST. Configuration **(BS-2CM)** applies MST to cluster regions and create text lines. Morphological text line generation is used in configuration **(BS-23M)**.

Modifications of Step (4): The following two configurations assess the methods for estimating the orientation of a text line: Configuration **(BS-4OP)** uses the Perpendicular Squared Distance method and configuration **(BS-4OS)** uses the Single String Orientation Detection method to estimate the orientation.

Modifications of Step (5): For all configurations, both OCR engines are used to generate the results. The identifier of a configuration is extended to **(BS-XYZ-T)** or **(BS-XYZ-O)**, when referencing the configurations that use Tesseract or Ocropy, respectively. Furthermore, we assess the direct impact of the OCR engine on the recognition results with configuration **(BS15-O)**, which only differs with respect to the OCR method from the base configuration by using the Ocropy OCR engine instead of Tesseract.

Modifications of Step (6): The last step of the pipeline is the post-processing. We use three configurations to evaluate the different post-processing methods: Configuration **(BS-6PC)** uses the Special Character Filter method for post-processing. Configuration **(BS-6PS)** uses the String Filter method for post-processing. Configuration **(BS-6PQ)** uses the Quantitative Assessment method for post-processing.

6 Evaluation

We conduct our evaluation on four datasets which are described in Section 6.1. The evaluation measures are described in Section 6.2.

Table 1. Number of figures, average figure width and height, and average number of text elements (TE), words, and characters per figure.

Dataset	# Figures	Width	Height	# TE	# Words	# Characters
EconBiz	121	982	681	25	35	151
DeGruyter	120	959	619	24	34	149
CHIME-R	115	714	454	14	18	69
CHIME-S	85	440	320	12	18	76
Total	441	801	535	19	27	114

6.1 Datasets

To the best of our knowledge, no high quality datasets for text extraction from scholarly figures exist. Thus, using our own tool that was specifically designed for manually labeling text elements in figures, we have created two datasets. One is in the domain of economics, the other is based on educational books. In addition, we use the CHIME datasets, which consist of figures of varying quality and origin. In total, we use four datasets with 441 figures of varying origin and characteristics.

- **EconBiz** We have created a corpus of 121 scholarly figures from the economics domain. We obtained these figures from a corpus of 288,000 open access publications from EconBiz¹⁰ by extracting all images, filtering them by size and other constraints, and randomly selecting the subset of 121 figures. The dataset resembles a wide variety of scholarly figures from bar charts to maps. The figures were manually labeled to create the necessary gold standard information.
- **DeGruyter** We manually labeled another dataset, composed of scholarly figures from books provided by DeGruyter¹¹ under a creative commons license¹². We selected ten books, mostly from the chemistry domain, which contain figures with English text and selected 120 figures randomly from these books. The gold standard for these figures was created using the same tool which has been used for the creation of the EconBiz dataset.
- **CHIME-R** The Chart Image Dataset¹³ consists of two subsets. The CHIME-R dataset consists of 115 real images that were collected on the Internet or scanned from paper. Most of the figures are bar charts. The rest are some pie charts and line charts. The gold standard was created by Yang Li [32].
- **CHIME-S** The other, CHIME-S dataset consists of 85 synthetically generated images. This set mainly contains line charts and pie charts and few bar charts. The gold standard was created by Zhao Jiuzhou [18].

¹⁰ <https://www.econbiz.de/>

¹¹ <http://www.degruyter.com/>

¹² <http://www.degruyter.com/dg/page/open-access-policy>

¹³ <https://www.comp.nus.edu.sg/~tancl/ChartImageDataset.htm>

Some statistics about the datasets which are useful for understanding the evaluation results can be found in Table 1. Both, the CHIME-R and CHIME-S datasets contain figures with an on average lower resolution than the EconBiz and DeGruyter datasets, which are almost equal. With respect to the average number of characters, words, and text elements in a figure, the distribution of EconBiz and DeGruyter are similar, with about twice as many as the CHIME datasets.

The use of standard test corpora greatly improves comparability. Thus, we looked at ImageNet, TREC, ICDAR, and ImageCLEF datasets, but none of these provide datasets for text extraction from scholarly figures. The closest candidates are the Medical Compound Figure Separation¹⁴ challenge from ImageCLEF and the image-to-structure task¹⁵ in the TREC chemistry track from 2011. The ImageCLEF dataset has the problem that the compound figures do not only contain scholarly figures but also diagnostic images like x-ray or ultrasound images. For TREC 2011, we focus on the chemistry track which aims at extracting the visually encoded structure of chemical elements. Thus, the chemistry task actually differs quite a lot from our general text extraction task. The text contained in the chemical images is even sparser and each string often consists of less than three characters. Furthermore, text and graphics are much more intertwined with a dominance in graphical components. The Robust Reading Competition (RRC) [19] held at the ICDAR conferences has a Born-Digital Images track which consists of artificial images extracted from the Web. We did not consider this dataset because none of its images are scholarly figures, therefore posing different challenges. For example, RRC images span from spam mails to logos and advertisements from websites. They contain only horizontal text, use specific unique fonts, with characters that are merged or broken by design. Furthermore, the images are often of low resolution. Extracting text at different orientations is a challenge common to scholarly figures that can not be evaluated using the ICDAR RRC dataset. This is also reflected in the 2013 report of the ICDAR RRC competition [20] where Karatzas et al. state that "All bounding boxes are defined as axis-aligned isothetic rectangles, a reasonable design choice as in born-digital images most text is horizontal.", which does not hold for scholarly figures. In addition, they leave out all words that have less than three characters. The ICDAR RRC evaluation consists of four tasks for measuring the performance. The first task is the evaluation of the text localization, which is based on the framework proposed by Wolf and Jolion [30]. The second task in the ICDAR RRC evaluation is text segmentation on the atom level as defined by Clavelli et al. [10]. The third evaluation task is word recognition, which works with a case-sensitive standard edit distance. The final task is a combination of the three previous tasks. The metrics used in Task 1, 2, and 4 are classic precision, recall, and F1-measure.

Based on the ICDAR RRC evaluation scheme, we have developed three measures, specifically designed for the structure of our gold standard, to evaluate our pipeline configurations.

¹⁴ <http://www.imageclef.org/2016/medical>

¹⁵ https://wiki.ir-facility.org/index.php/Image2Structure_Task

6.2 Measures

We have selected four measures to evaluate the pipeline configurations and compare their results. Our gold standard consists of text elements which represent single lines of text taken from a scholarly figure. Each text line consists of one or multiple words which are separated by blank space. Each word may consist of any combination of characters and numbers. Every text line is defined by a specific position, size, and orientation.

Each pipeline configuration generates a set of text line elements as well. These text lines need to be matched to the gold standard. Since we do not have pixel information per character, we match the extraction results with the gold standard by using the bounding boxes. This is based on the first evaluation task of the ICDAR RRC and evaluates the text localization on text line level. We iterate over all text lines in the gold standard and take all matches that are above the so-called intersection threshold. Our matching procedure calculates the intersection area between all pairs of the pipeline output and gold standard text lines. If the intersection comprises at least ten percent of the combined area of both text elements, then it is considered a match. This reduces the error introduced through elements which are an incorrect match and only have a small overlap with the gold standard, but still tolerates it when a line was broken up into multiple parts. We look at each gold standard element and take all elements from the pipeline as matches that are above the intersection threshold. Thus, a gold standard element can have multiple matching elements and an element from the pipeline can be assigned to multiple elements from the gold standard if it fulfills the matching constraint for each match. We have defined three measures to assess these matches. The first two measures analyze the text localization. The third measure compares the recognized text, similar to the word recognition task of the ICDAR RRC, although we compare text lines and not individual words. Finally, we also decided to analyze the performance of the individual methods.

Text Location Detection First, we evaluate how accurate the configurations are at detecting the text locations. If at least one match is found for an element from the gold standard set, it counts as a true positive, regardless of what text was recognized. If no match was found, it is considered as false negative. A false positive is an element from the pipeline output which has no match. From these values, we compute precision, recall, and F1-measure for assessing the text location detection. This measure is a binary evaluation and assesses only whether a match to an element exists or not. Additionally, we report the Element Ratio (ER) which is the number of elements recognized by the pipeline divided by the number of elements in the gold standard and the Matched Element Ratio (MER) which is the number of matched items from the pipeline divided by the number of elements of the gold standard. These ratios give an idea whether gold standard elements get matched by multiple elements and whether the configuration tends to find more elements or less elements than it actually should find. We restrict the number of matches by applying a coverage threshold, which means that the intersection area of two elements has to be at least as big as 10% of the

total covered area. This measure penalizes mappings between elements of largely varying size. We further evaluate these matches with the next two measures.

Text Element Coverage Second, we investigate the matching in more detail by assessing the text element coverage. For each gold standard text element, we take the pixels of the bounding boxes and compute their overlap to calculate precision, recall, and F1-measure over all of its matches. The true positives in this case are the overlapping pixels and the false positives are those pixels from the text elements from the pipeline which are not overlapping. The false negatives are the pixels of the gold standard element which were not covered by a text element from the pipeline. The values are averaged over all gold standard text elements in a figure.

Text Recognition Quality Third, we assess the quality of the recognized text by computing the Levenshtein distance between the extracted text and the gold standard. We calculate the distance for each match and report the average for the whole figure. Since multiple text elements from the pipeline can be matched to a gold standard text line, we have to combine their text into one string. We combine the elements using their position information. Besides a (local) Levenshtein Distance per match, we also compute a global Levenshtein distance over all extracted text. This means that for each figure, we combine all characters from the text elements of the gold standard and add them to one string. Likewise, we create a string from the text elements extracted by the pipeline. The characters in both strings are sorted alphabetically and we compute the Levenshtein Distance between these strings. This approximates the overall number of operations needed to match the strings without considering position information. Since the global Levenshtein Distance depends on the number of characters inside a figure, we also report an operations per character (OPC) score, which is computed by dividing the global Levenshtein Distance by the number of characters in the gold standard. This normalizes the global Levenshtein Distance and makes it comparable across scholarly figures with different amounts of characters.

Runtime Performance Finally, we log the execution time for each method in milliseconds to see which method has the best runtime performance. We analyze each configuration individually and compute the average time for each method over all figures from all four datasets. For each configuration, we aggregate the execution time of its methods to compute the average time needed per step of the pipeline.

7 Results

We have executed all configurations listed in Section 5 and analyzed the output with respect to the measures for Text Location Detection, Text Element Coverage, Text Recognition Quality, and Runtime Performance as defined in Section 6.2. For reasons of simplicity, we are only reporting the average values

over all datasets for all configurations as well as for each dataset separately. We compute the average Precision/Recall/F1-measure over the elements of each figure. We report the average Precision/Recall/F1-measure in terms of mean and standard deviation over all individual results per figure. The local Levenshtein distance is reported as the average of the mean values per figure and the average standard deviation. The global Levenshtein distance is defined by the mean and standard deviation over all figures and the normalized OPC score. Subsequently, the performance measures are reported. First, we report the results of the configurations from the literature. Subsequently, we present the results for the systematically modified configurations.

Text Location Detection, Text Element Coverage, and Text Recognition Quality The text location detection results for the configurations from the literature computed over all datasets are reported in Table 2. The best result, based on the F1-measure, is achieved by configuration **(BS15)** with a F1-measure of 0.58. Looking at each dataset separately, as documented in Table 3, one can see that configuration **(BS15)** works best for the DeGruyter (0.70) and CHIME-R (0.63) datasets while **(SZ13)** has the best result for the CHIME-S (0.55) and EconBiz (0.57) datasets. The coverage assessment in Table 4 shows the best precision of 0.79 for **(Hu05)**, the best recall of 0.59 for **(SZ13)**, and the best F1-measure of 0.57 for **(Hu05)**. The individual results per dataset for all configurations are shown in Table 5. The text recognition quality is presented in Table 6 and the individual results per dataset for all configurations are presented in Table 7. We obtain the best results with **(BS15)** with only 0.67 operations per character (OPC), an average local Levenshtein of 6.23, and an average global Levenshtein of 108.81. Only configuration **(CK15)** has a slightly better average local Levenshtein (6.07).

Table 2. Configurations from the literature: Precision (Pr), Recall (Re), and F1-measure for the Average Text Location Detection, Element Ratio (ER), and Matched Element Ratio (MER). Results are averaged over all datasets.

Config.	Pr	Re	F1 (SD)	ER	MER
SZ13	0.63	0.47	0.54 (0.23)	0.80	0.59
Hu05	0.61	0.43	0.48 (0.28)	0.77	0.57
Ja07	0.59	0.45	0.49 (0.28)	0.83	0.51
BS15	0.66	0.55	0.58 (0.25)	1.04	0.69
CK15	0.52	0.50	0.53 (0.23)	1.37	0.60
Fr15	0.55	0.51	0.54 (0.25)	1.44	0.72
XK10	0.73	0.35	0.45 (0.26)	0.43	0.39

For the systematically modified configurations, Table 8 shows the text location detection results, Table 9 the coverage assessment, and Table 10 shows the text recognition results. The best location detection F1-measure of 0.67 is

achieved by **(BS-4OS)**, which is also supported by the coverage assessment in Table 9 with the highest F1-measure of 0.65. The separate evaluation of each dataset in Table 3 confirms as well that **(BS-4OS)** works best for EconBiz, CHIME-R, and CHIME-S with F1-measures of 0.68, 0.66, and 0.63. The best result for DeGruyter is an F1-measure of 0.73 by **(BS-23M)** with **(BS-4OS)** having the second best F1-measure of 0.71. The best coverage assessment results for EconBiz (0.56), DeGruyter (0.67), CHIME-R (0.71), and CHIME-S (0.66) are by **(BS-4OS)**, as shown in Table 5. Configuration **(BS-4OS-O)** also pro-

Table 3. Average F1-measure and Standard Deviation for Text Location Detection per Configuration on each Dataset

Config.	EconBiz	DeGruyter	CHIME-R	CHIME-S
SZ13	0.57(0.25)	0.51(0.24)	0.53(0.21)	0.55(0.23)
Hu05	0.45(0.29)	0.55(0.30)	0.50(0.25)	0.34(0.26)
Ja07	0.47(0.29)	0.50(0.27)	0.52(0.26)	0.33(0.22)
CK15	0.53(0.22)	0.54(0.21)	0.56(0.24)	0.49(0.25)
Fr15	0.48(0.26)	0.62(0.21)	0.58(0.24)	0.41(0.21)
XK10	0.38(0.24)	0.50(0.24)	0.48(0.26)	0.29(0.18)
BS15	0.55(0.25)	0.70(0.18)	0.63(0.23)	0.43(0.25)
BS-1NC	0.54(0.26)	0.65(0.20)	0.61(0.24)	0.42(0.26)
BS-1OC	0.46(0.28)	0.51(0.28)	0.52(0.24)	0.43(0.25)
BS-1QP	0.42(0.23)	0.53(0.23)	0.49(0.27)	0.36(0.22)
BS-2nF	0.46(0.22)	0.59(0.19)	0.53(0.22)	0.37(0.25)
BS-2CG	0.48(0.29)	0.64(0.22)	0.58(0.27)	0.33(0.28)
BS-2CM	0.55(0.25)	0.70(0.21)	0.62(0.22)	0.40(0.24)
BS-23M	0.62(0.24)	0.73(0.17)	0.63(0.22)	0.47(0.25)
BS-4OP	0.57(0.24)	0.64(0.20)	0.59(0.24)	0.39(0.27)
BS-4OS	0.68(0.21)	0.71(0.18)	0.66(0.23)	0.63(0.26)
BS-6PC	0.55(0.25)	0.71(0.17)	0.62(0.23)	0.43(0.26)
BS-6PS	0.55(0.26)	0.71(0.17)	0.64(0.23)	0.43(0.25)
BS-6PQ	0.38(0.24)	0.57(0.21)	0.50(0.21)	0.36(0.25)

Table 4. Configurations motivated the literature: Precision (Pr), Recall (Re), and F1-measure for the Average Text Element Coverage. Results are averaged over all datasets.

Config.	<i>Pr</i>	<i>Re</i>	<i>F1 (SD)</i>
SZ13	0.52	0.59	0.47 (0.21)
Hu05	0.79	0.54	0.57 (0.20)
Ja07	0.41	0.32	0.32 (0.21)
BS15	0.60	0.49	0.50 (0.24)
CK15	0.53	0.41	0.42 (0.21)
Fr15	0.65	0.54	0.54 (0.23)
XK10	0.33	0.34	0.30 (0.22)

duces the best text recognition results with an average local Levenshtein of 4.71 and an OPC of 0.53. The best local Levenshtein for each dataset is also achieved by configuration **(BS-4OS-O)** with values between 3.51 and 5.80 (see Table 7). In addition, configuration **(BS-4OS-O)** shows the best results of 95.49 for the average global Levenshtein Distance. Comparing the different, systematically modified configurations per step of the pipeline shows that the only major improvement is achieved by **(BS-4OS)**.

Table 5. Average F1-measure and Standard Deviation for Text Element Coverage per Configuration on each Dataset

Config.	EconBiz	DeGruyter	CHIME-R	CHIME-S
SZ13	0.42(0.18)	0.44(0.23)	0.49(0.21)	0.55(0.20)
Hu05	0.48(0.17)	0.58(0.22)	0.66(0.18)	0.47(0.22)
Ja07	0.28(0.19)	0.33(0.20)	0.41(0.22)	0.21(0.19)
CK15	0.37(0.19)	0.48(0.20)	0.45(0.22)	0.35(0.23)
Fr15	0.42(0.21)	0.62(0.15)	0.62(0.25)	0.45(0.25)
XK10	0.23(0.17)	0.34(0.21)	0.37(0.26)	0.16(0.11)
BS15	0.40(0.20)	0.62(0.14)	0.60(0.26)	0.31(0.21)
BS-1NC	0.36(0.21)	0.58(0.16)	0.58(0.26)	0.27(0.17)
BS-1OC	0.32(0.24)	0.46(0.28)	0.42(0.25)	0.31(0.21)
BS-1QP	0.43(0.22)	0.45(0.23)	0.43(0.26)	0.33(0.22)
BS-2nF	0.40(0.17)	0.56(0.16)	0.60(0.22)	0.35(0.23)
BS-2CG	0.46(0.19)	0.58(0.16)	0.69(0.19)	0.45(0.23)
BS-2CM	0.36(0.20)	0.60(0.16)	0.59(0.25)	0.24(0.18)
BS-23M	0.39(0.19)	0.62(0.14)	0.56(0.22)	0.30(0.17)
BS-4OP	0.37(0.17)	0.47(0.14)	0.48(0.21)	0.20(0.20)
BS-4OS	0.56(0.14)	0.67(0.11)	0.71(0.21)	0.66(0.20)
BS-6PC	0.39(0.20)	0.62(0.15)	0.60(0.26)	0.30(0.20)
BS-6PS	0.39(0.20)	0.61(0.15)	0.59(0.26)	0.31(0.21)
BS-6PQ	0.21(0.16)	0.41(0.20)	0.36(0.22)	0.19(0.17)

Table 6. Average local Levenshtein (L) and global Levenshtein (G) and Operations Per Character (OPC) over all datasets for the configurations from the literature using Tesseract

Config.	$AVG_L(SD)$	$AVG_G(SD)$	OPC
SZ13	6.67 (4.82)	122.28 (141.03)	0.70
Hu05	6.65 (5.41)	126.35 (138.95)	0.71
Ja07	7.92 (5.56)	150.25 (140.59)	1.13
BS15	6.23 (4.93)	108.81 (108.53)	0.67
CK15	6.07 (5.08)	120.12 (125.87)	0.71
Fr15	6.72 (6.02)	135.64 (201.31)	0.85
XK10	7.06 (5.41)	125.45 (134.88)	0.74

Runtime Performance Figure 4 compares the seven configurations motivated from the literature. For each configuration, the average execution time over all figures from all four datasets are reported. The execution time is split up along the pipeline steps. Please note that the runtime is specified in milliseconds. Figure 4 shows the results at log scale, due to the varying performance. As one can see from the figure, the configurations (**SZ13**), (**Hu05**), and (**Ja07**) are in the same range, with the latter one being the fastest configuration. The other configurations are by an order of magnitude slower, with (**CK15**) being the configuration needing most of the time (about 38 days for all four datasets). One can also clearly see that in general the first step of the pipeline contributes most to the execution time. Looking at the data generated by the systematically modified configurations, we make the following observations. The fastest binarization is Otsu’s method with on average less than 100ms and adaptive binarization requires the most time with about two minutes on average. The pivoting algorithm requires about twice as much time as Connected Component Labeling. Looking at step 2 and 3 of the pipeline, only the morphological clustering differs a lot from the rest with an average execution time of several minutes or more. The PSD and Hough orientation estimation execute on average in a few milliseconds, while the SSOD needs several seconds. The comparison of Ocropy and Tesseract

Table 7. Average Levenshtein and Standard Deviation for Text Recognition Quality per Configuration on each Dataset

Config.	EconBiz		DeGruyter		CHIME-R		CHIME-S	
	T	O	T	O	T	O	T	O
SZ13	6.00(3.24)	-	6.13(3.25)	-	7.10(6.25)	-	7.29(5.64)	-
Hu05	5.69(3.33)	-	5.69(3.74)	-	6.75(6.95)	-	7.94(6.13)	-
Ja07	7.15(3.62)	-	7.97(3.85)	-	7.94(7.25)	-	8.36(6.34)	-
CK15	5.27(3.11)	-	4.96(2.78)	-	6.65(7.21)	-	6.74(5.30)	-
Fr15	5.74(3.34)	-	6.53(5.87)	-	6.81(7.76)	-	7.26(5.55)	-
XK10	6.27(3.46)	-	6.46(3.56)	-	6.99(7.14)	-	7.88(6.09)	-
BS15	5.42(3.06)	4.80(2.93)	4.88(2.58)	4.07(2.26)	6.51(6.85)	5.96(6.97)	7.21(5.34)	7.20(5.66)
BS-1NC	5.47(3.12)	4.81(2.86)	5.30(2.78)	4.47(2.83)	6.42(6.85)	5.94(7.00)	7.22(5.42)	7.47(5.71)
BS-1OC	5.74(3.20)	5.16(3.10)	5.76(3.05)	5.21(3.23)	6.72(6.72)	6.54(6.89)	7.51(5.71)	7.35(5.89)
BS-1QP	7.13(4.13)	6.18(3.33)	8.44(4.88)	6.29(3.04)	8.50(7.88)	7.19(7.72)	8.68(6.49)	8.51(6.27)
BS-2nF	5.78(3.33)	5.30(2.94)	5.83(3.11)	5.02(2.58)	6.64(6.29)	6.91(7.83)	7.63(5.72)	7.78(6.22)
BS-2CG	5.77(3.43)	5.28(3.27)	5.49(3.79)	4.71(2.92)	6.63(7.46)	6.31(7.83)	8.34(6.27)	8.14(6.29)
BS-2CM	5.57(3.07)	4.98(3.00)	4.95(2.75)	4.45(2.42)	6.72(7.27)	6.09(7.35)	7.51(6.08)	7.73(6.10)
BS-23M	5.17(3.16)	4.56(2.89)	4.92(3.34)	3.82(2.24)	6.66(6.68)	6.13(7.01)	7.53(5.62)	7.33(5.60)
BS-4OP	7.56(3.95)	6.34(3.29)	8.82(3.84)	6.71(3.44)	8.15(7.25)	7.36(7.81)	8.49(6.47)	8.22(6.44)
BS-4OS	4.90(3.01)	3.86(2.82)	4.55(2.55)	3.51(2.00)	5.76(6.01)	5.08(6.32)	6.27(4.97)	5.80(5.50)
BS-6PC	5.13(3.03)	4.73(2.93)	4.77(2.33)	4.28(2.23)	6.06(6.71)	5.73(6.82)	7.29(5.43)	7.19(5.63)
BS-6PS	5.39(3.02)	4.75(2.89)	4.83(2.44)	4.08(2.27)	6.53(6.87)	5.79(6.87)	7.24(5.37)	7.16(5.61)
BS-6PQ	5.47(3.13)	5.20(2.99)	4.96(2.63)	4.64(2.77)	6.31(6.91)	5.95(6.59)	7.43(5.70)	7.57(6.05)

Table 8. Systematically modified configurations: Precision (Pr), Recall (Re), and F1-measure for the Average Text Location Detection, Element Ratio (ER), and Matched Element Ratio (MER). Results are averaged over all datasets

Config.	<i>Pr</i>	<i>Re</i>	<i>F1 (SD)</i>	<i>ER</i>	<i>MER</i>
BS15	0.66	0.55	0.58 (0.25)	1.04	0.69
BS-1NC	0.64	0.52	0.57 (0.25)	0.96	0.64
BS-1OC	0.67	0.40	0.49 (0.26)	0.74	0.53
BS-1QP	0.61	0.44	0.48 (0.25)	0.96	0.75
BS-2nF	0.60	0.46	0.51 (0.23)	0.86	0.52
BS-2CG	0.62	0.50	0.55 (0.27)	0.90	0.64
BS-2CM	0.61	0.54	0.59 (0.25)	1.19	0.74
BS-23M	0.67	0.55	0.62 (0.23)	1.08	0.65
BS-4OP	0.62	0.53	0.57 (0.24)	1.01	0.66
BS-4OS	0.67	0.63	0.67 (0.22)	1.27	0.88
BS-6PC	0.69	0.54	0.59 (0.25)	0.97	0.70
BS-6PS	0.67	0.55	0.60 (0.25)	1.01	0.69
BS-6PQ	0.66	0.38	0.48 (0.25)	0.60	0.43

shows that the latter one is about three times faster. Finally, all post-processing methods need on average less than a millisecond.

Table 9. Systematically modified configurations: Precision (Pr), Recall (Re), and F1-measure for the Average Text Element Coverage. Results are averaged over all datasets.

Config.	<i>Pr</i>	<i>Re</i>	<i>F1 (SD)</i>
BS15	0.60	0.49	0.50 (0.24)
BS-1NC	0.59	0.44	0.47 (0.24)
BS-1OC	0.46	0.40	0.38 (0.26)
BS-1QP	0.41	0.57	0.42 (0.23)
BS-2nF	0.59	0.54	0.50 (0.21)
BS-2CG	0.76	0.54	0.57 (0.20)
BS-2CM	0.57	0.47	0.47 (0.24)
BS-23M	0.60	0.47	0.48 (0.22)
BS-4OP	0.49	0.40	0.41 (0.20)
BS-4OS	0.77	0.63	0.65 (0.17)
BS-6PC	0.59	0.49	0.49 (0.24)
BS-6PS	0.59	0.49	0.49 (0.24)
BS-6PQ	0.39	0.29	0.31 (0.21)

Table 10. Average local Levenshtein (L) and global Levenshtein (G) and Operations Per Character (OPC) over all datasets for the systematic configurations

Config.	Tesseract			Ocropy		
	$AVG_L(SD)$	$AVG_G(SD)$	OPC	$AVG_L(SD)$	$AVG_G(SD)$	OPC
BS15	6.23 (4.93)	108.81 (108.53)	0.67	5.47 (4.98)	108.55 (106.64)	0.64
BS-1NC	6.27 (4.95)	117.58 (124.23)	0.69	5.70 (5.09)	117.46 (128.73)	0.66
BS-1OC	6.55 (5.06)	131.58 (142.74)	0.75	6.16 (5.21)	131.39 (143.16)	0.73
BS-1QP	8.31 (6.14)	154.54 (168.10)	1.09	7.06 (5.62)	136.40 (132.05)	0.82
BS-2nF	6.55 (4.94)	111.30 (105.13)	0.75	6.29 (5.50)	120.71 (109.18)	0.76
BS-2CG	6.68 (5.65)	108.86 (102.93)	0.66	6.22 (5.75)	130.21 (127.87)	0.69
BS-2CM	6.30 (5.29)	115.43 (113.79)	0.69	5.85 (5.34)	110.74 (107.23)	0.67
BS-23M	6.15 (5.12)	104.61 (105.97)	0.63	5.52 (5.10)	106.71 (104.05)	0.64
BS-4OP	8.30 (5.59)	147.91 (129.55)	1.04	7.23 (5.60)	135.21 (122.48)	0.85
BS-4OS	5.47 (4.39)	96.29 (99.44)	0.58	4.71 (4.66)	95.49 (94.80)	0.53
BS-6PC	5.96 (4.88)	105.50 (107.16)	0.61	5.46 (5.00)	109.07 (104.57)	0.63
BS-6PS	6.20 (4.90)	108.06 (109.38)	0.64	5.45 (4.96)	106.38 (103.29)	0.63
BS-6PQ	6.07 (5.03)	120.78 (122.44)	0.67	5.79 (4.97)	126.92 (124.06)	0.71

8 Discussion

Comparing the different configurations from the literature shows that the best performing configuration is **(BS15)**. A possible reason is that our pipeline does not make many assumptions about the figures. Thus performing better on the heterogeneous datasets. In the following, we will discuss the results for the individual pipeline steps based on the results from the systematically modified configurations. Comparing the configurations for the first pipeline step leads to the conclusion that the adaptive binarization works best, because it can adapt to local variations of the appearance in a figure. Otsu’s method is too simple and Niblack’s method is more suited for document images which have fewer color variations. The lower results for the pivoting algorithm can be explained with the larger regions and the possibility that a region can be a mixture of text and graphic elements due to the only horizontal and vertical subdivision. Looking at step 2 and 3 of the pipeline, only the morphological clustering shows slightly better results than the DBSCAN-MST combination, most likely due to its processing on pixel level. However, it is by an order of magnitude slower. The overall best results are achieved by **(BS-4OS)**. This can be explained by the fact that the orientation estimation via Hough in the base configuration works on the centers of mass of character regions, which is an aggregated region representation, while the SSOD in **(BS-4OS)** computes the orientation on the original pixels. Thus, it avoids a possible error, induced by the pixel aggregation, which comes at the price of a higher computational time like with the morphological clustering. When comparing the OCR engines from step 5, Ocropy generally produces better results than Tesseract. Ocropy seems to be more conservative, having built in much more restrictions about what input to accept and when to execute

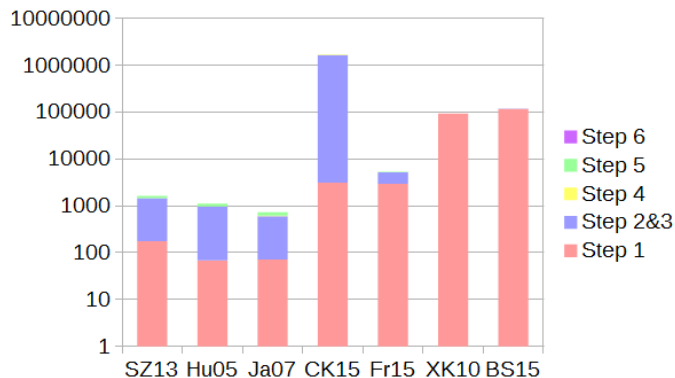


Fig. 4. Average execution time aggregated per pipeline step for the configurations motivated from the related work.

the OCR. Furthermore, each OCR engine comes with its own English language model and we did not evaluate their influence. The methods for post-processing do not seem to be an improvement. One reason might be the simplicity of their nature. Thus, some more advanced techniques may be developed in the future.

The performance analysis results for the configurations motivated from the literature clearly show that the first pipeline step is the most expensive one with respect to the execution time. Those configurations which use standard Otsu’s method are a lot faster than those using Adaptive Otsu Binarization. This difference was expected, since the latter one recursively processes a figure while applying Otsu’s method multiple times. In addition, configuration (**CK15**) did not finish in reasonable time (about 2 hours per image, 38 days in total) due to its expensive morphological clustering. This method is very expensive since it iterates multiple times over an image increasing the region by single pixel borders each time. The systematic comparisons confirm these results. Furthermore, the systematic comparisons show that the SSOD is more expensive than Hough and PSD, because it works on the pixel level similar to the morphological clustering. Furthermore, the difference between Ocropy and Tesseract can be explained by the fact that Tesseract is used via an API while Ocropy is currently used via command line which requires system calls. The fast execution of the post-processing methods results from their simplicity, applying only simple string comparisons.

Overall, there are many more options for the different pipeline steps, e. g., other binarization methods, different clustering algorithms, or post-processing methods. However, we made a selection of relevant approaches and methods to limit the combinatorial complexity. Furthermore, we did not consider commercial OCR engines like AbbyyFineReader or OmnipageReader, nor GOCR (which is only rarely updated, last update from 2013). Since adding more methods would lead to additional pipeline configurations, which significantly increases

the complexity of our experiments. However, as stated in the introduction, we provide the datasets and the implementation of the generic pipeline that was used in our experiment to the public. This allows for integrating and comparing new methods as well as the reproduction of our results.

9 Conclusion

Our comparison of the 32 pipeline configurations for text extraction shows that there is a clear favorite, configuration **(BS-4OS)**. The systematic modifications revealed that the SSOD line orientation estimation works better than the Hough Transformation and Ocropus returns better OCR results than Tesseract, when comparing with the best configuration from the literature. However, there is still room for other improvements: We plan to further expand our experiment by integrating commercial OCR engines like AbbyyFineReader. Furthermore, the use of heuristic post-processing has only little influence on the final results and other methods should be investigated. In addition, further extensions with respect to superscripts and subscripts as well as mathematical formulas should be investigated. Finally, we also plan to investigate non-linear pipeline configurations and extend our pipeline to target the ICDAR Robust Reading Competition on the Born-Digital Images dataset.

Acknowledgement This research was co-financed by the EU H2020 project MOVING (<http://www.moving-project.eu/>) under contract no 693092.

References

1. F. Bösch and A. Scherp. Formalization and preliminary evaluation of a pipeline for text extraction from infographics. In R. Bergmann, S. Görg, and G. Müller, editors, *Proceedings of the LWA 2015 Workshops: KDML, FGWM, IR, and FGDB, Trier, Germany, October 7-9, 2015.*, volume 1458 of *CEUR Workshop Proceedings*, pages 20–31. CEUR-WS.org, 2015.
2. F. Bösch and A. Scherp. Multi-oriented text extraction from information graphics. In C. Vanoirbeek and P. Genevès, editors, *Proceedings of the 2015 ACM Symposium on Document Engineering, DocEng 2015, Lausanne, Switzerland, September 8-11, 2015*, pages 35–38. ACM, 2015.
3. S. Carberry, S. Elzer, and S. Demir. Information graphics: an untapped resource for digital libraries. In E. N. Efthimiadis, S. T. Dumais, D. Hawking, and K. Järvelin, editors, *SIGIR 2006: Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Seattle, Washington, USA, August 6-11, 2006*, pages 581–588. ACM, 2006.
4. S. Carberry, S. E. Schwartz, K. F. McCoy, S. Demir, P. Wu, C. F. Greenbacker, D. Chester, E. Schwartz, D. Oliver, and P. S. Moraes. Access to multimodal articles for individuals with sight impairments. *TiiS*, 2(4):21, 2012.
5. Z. Chen, M. J. Cafarella, and E. Adar. Diagramflyer: A search engine for data-driven diagrams. In A. Gangemi, S. Leonardi, and A. Panconesi, editors, *Proceedings of the 24th International Conference on World Wide Web Companion, WWW 2015, Florence, Italy, May 18-22, 2015 - Companion Volume*, pages 183–186. ACM, 2015.

6. D. Chester and S. Elzer. Getting computers to see information graphics so users do not have to. In M. Hacid, N. V. Murray, Z. W. Ras, and S. Tsumoto, editors, *Foundations of Intelligent Systems, 15th International Symposium, ISMIS 2005, Saratoga Springs, NY, USA, May 25-28, 2005, Proceedings*, volume 3488 of *Lecture Notes in Computer Science*, pages 660–668. Springer, 2005.
7. Y. Chiang and C. A. Knoblock. A general approach for extracting road vector data from raster maps. *IJDAR*, 16(1):55–81, 2013.
8. Y. Chiang and C. A. Knoblock. Recognizing text in raster maps. *GeoInformatica*, 19(1):1–27, 2015.
9. S. R. Choudhury and C. L. Giles. An architecture for information extraction from figures in digital libraries. In A. Gangemi, S. Leonardi, and A. Panconesi, editors, *Proceedings of the 24th International Conference on World Wide Web Companion, WWW 2015, Florence, Italy, May 18-22, 2015 - Companion Volume*, pages 667–672. ACM, 2015.
10. A. Clavelli, D. Karatzas, and J. Lladós. A framework for the assessment of text extraction algorithms on complex colour images. In *Document Analysis Systems*, ACM International Conference Proceeding Series, pages 19–26. ACM, 2010.
11. M. P. Deseilligny, H. L. Men, and G. Stamon. Character string recognition on maps, a rotation-invariant recognition method. *Pattern Recognition Letters*, 16(12):1297–1310, 1995.
12. M. Fraz, M. S. Sarfraz, and E. A. Edirisinghe. Exploiting colour information for better scene text detection and recognition. *IJDAR*, 18(2):153–167, 2015.
13. J. Gllavata and B. Freisleben. Adaptive fuzzy text segmentation in images with complex backgrounds using color and texture. In A. Gagalowicz and W. Philips, editors, *Computer Analysis of Images and Patterns, 11th International Conference, CAIP 2005, Versailles, France, September 5-8, 2005, Proceedings*, volume 3691 of *Lecture Notes in Computer Science*, pages 756–765. Springer, 2005.
14. W. Huang and C. L. Tan. A system for understanding imaged infographics and its applications. In P. R. King and S. J. Simske, editors, *Proceedings of the 2007 ACM Symposium on Document Engineering, Winnipeg, Manitoba, Canada, August 28-31, 2007*, pages 9–18. ACM, 2007.
15. W. Huang, C. L. Tan, and W. K. Leow. Associating text and graphics for scientific chart understanding. In *Eighth International Conference on Document Analysis and Recognition (ICDAR 2005), 29 August - 1 September 2005, Seoul, Korea*, pages 580–584. IEEE Computer Society, 2005.
16. J. Illingworth and J. Kittler. A survey of the hough transform. *Computer Vision, Graphics, and Image Processing*, 44(1):87–116, 1988.
17. C. Jayant, M. Renzelmann, D. Wen, S. Krisnandi, R. E. Ladner, and D. Comden. Automated tactile graphics translation: in the field. In E. Pontelli and S. Trewin, editors, *Proceedings of the 9th International ACM SIGACCESS Conference on Computers and Accessibility, ASSETS 2007, Tempe, Arizona, USA, October 15-17, 2007*, pages 75–82. ACM, 2007.
18. Z. Jiuzhou. Creation of synthetic chart image database with ground truth. Honors year project report, National University of Singapore, 2006. https://www.comp.nus.edu.sg/~tancl/ChartImageDatabase/Report_Zhaojiuzhou.pdf.
19. D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, S. K. Ghosh, A. D. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu, F. Shafait, S. Uchida, and E. Valveny. ICDAR 2015 competition on robust reading. In *13th International Conference on Document Analysis and Recognition, ICDAR 2015, Nancy, France, August 23-26, 2015*, pages 1156–1160. IEEE Computer Society, 2015.

20. D. Karatzas, F. Shafait, S. Uchida, M. Iwamura, L. G. i Bigorda, S. R. Mestre, J. Mas, D. F. Mota, J. Almazán, and L. de las Heras. ICDAR 2013 robust reading competition. In *ICDAR*, pages 1484–1493. IEEE Computer Society, 2013.
21. K. Khurshid, I. Siddiqi, C. Faure, and N. Vincent. Comparison of Niblack inspired binarization methods for ancient documents. In K. Berkner and L. Likforman-Sulem, editors, *Document Recognition and Retrieval XVI, DRR 2009, 16th Document Recognition and Retrieval Conference, part of the IS&T-SPIE Electronic Imaging Symposium, San Jose, CA, USA, January 18-22, 2009. Proceedings*, volume 7247 of *SPIE Proceedings*, pages 1–10. SPIE, 2009.
22. S. Lu, T. Chen, S. Tian, J. Lim, and C. L. Tan. Scene text extraction based on edges and support vector regression. *IJDAR*, 18(2):125–135, 2015.
23. X. Lu, S. Kataria, W. J. Brouwer, J. Z. Wang, P. Mitra, and C. L. Giles. Automated analysis of images in documents for intelligent document search. *IJDAR*, 12(2):65–81, 2009.
24. J. I. Olszewska. Active contour based optical character recognition for automated scene understanding. *Neurocomputing*, 161:65–71, 2015.
25. N. Otsu. A threshold selection method from gray-level histograms. *Systems, Man and Cybernetics, IEEE Transactions on*, 9(1):62–66, Jan 1979.
26. H. Samet and M. Tamminen. Efficient component labeling of images of arbitrary dimension represented by linear bintrees. *IEEE TPAMI*, 10(4):579–586, 1988.
27. J. Sas and A. Zolnierek. Three-stage method of text region extraction from diagram raster images. In R. Burduk, K. Jackowski, M. Kurzynski, M. Wozniak, and A. Zolnierek, editors, *Proceedings of the 8th International Conference on Computer Recognition Systems CORES 2013, Milkow, Poland, 27-29 May 2013*, volume 226 of *Advances in Intelligent Systems and Computing*, pages 527–538. Springer, 2013.
28. M. Savva, N. Kong, A. Chhajta, F. Li, M. Agrawala, and J. Heer. Revision: automated classification, analysis and redesign of chart images. In J. S. Pierce, M. Agrawala, and S. R. Klemmer, editors, *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology, Santa Barbara, CA, USA, October 16-19, 2011*, pages 393–402. ACM, 2011.
29. C. M. Strohmaier, C. Ringlstetter, K. U. Schulz, and S. Mihov. Lexical post-correction of ocr-results: The web as a dynamic secondary dictionary? In *7th International Conference on Document Analysis and Recognition (ICDAR 2003), 2-Volume Set, 3-6 August 2003, Edinburgh, Scotland, UK*, pages 1133–1137. IEEE Computer Society, 2003.
30. C. Wolf and J. Jolion. Object count/area graphs for the evaluation of object detection and segmentation algorithms. *IJDAR*, 8(4):280–296, 2006.
31. S. Xu and M. Krauthammer. A new pivoting and iterative text detection algorithm for biomedical images. *Journal of Biomedical Informatics*, 43:924–931, 2010.
32. L. Yang, W. Huang, and C. L. Tan. Semi-automatic ground truth generation for chart image recognition. In H. Bunke and A. L. Spitz, editors, *Document Analysis Systems VII, 7th International Workshop, DAS 2006, Nelson, New Zealand, February 13-15, 2006, Proceedings*, volume 3872 of *Lecture Notes in Computer Science*, pages 324–335. Springer, 2006.