



# Enhancing OCR in historical documents with complex layouts through machine learning

David Fleischhacker<sup>1</sup> · Roman Kern<sup>1,2</sup> · Wolfgang Göderle<sup>3,4</sup>

Received: 5 October 2023 / Revised: 20 August 2024 / Accepted: 19 January 2025  
© The Author(s) 2025

## Abstract

This paper explores the challenge of processing and extracting information from large quantities of printed serial sources from the 19th century, which have been largely untapped due to the inadequacies of existing extraction techniques. We focus on the Habsburg Central Europe's *Hof- und Staatsschematismus*, a comprehensive record published between 1702 and 1918 that documents the Habsburg civil service's hierarchy and the evolution of its central administration over two centuries. Our approach sees the significant investment into machine learning-driven layout detection prior to the OCR-process. We generated synthetic data mimicking the *Hof- und Staatsschematismus* style for initial training of a Faster R-CNN model, followed by fine-tuning the model with a smaller dataset of manually annotated historical documents. Subsequently, we optimised Tesseract-OCR for our document style to enhance the combined structure extraction and OCR process. Our evaluation demonstrates significant improvements in OCR performance metrics (WER and CER), with the combined structure detection and fine-tuned OCR process showing a decrease in error rates of 15.68 percentage points for CER and 19.95 percentage points for WER. These findings underscore the potential of ML techniques in facilitating the extraction and analysis of historical documents.

**Keywords** PDF extraction · Layout detection · OCR fine-tuning · Synthetic training data · Document analysis and recognition

## 1 Introduction

The long 19th century provides historians and fellow humanists with a wealth of retrodigitized printed sources. A significant share of these is made up of serial publications

---

Roman Kern and Wolfgang Göderle have contributed equally to this work.

---

✉ David Fleischhacker  
david.fleischhacker@student.tugraz.at

Roman Kern  
rkern@tugraz.at

Wolfgang Göderle  
wolfgang.goederle@uibk.ac.at

<sup>1</sup> Institute of Interactive Systems and Data Science, Graz University of Technology, Sandgasse 36, 8010 Graz, Austria

<sup>2</sup> Know Center, Sandgasse 36, 8010 Graz, Austria

<sup>3</sup> Department of Austrian History, Institute of History, University of Innsbruck, Innrain 52, 6020 Innsbruck, Austria

<sup>4</sup> Department Structural Changes of the Technosphere, Max Planck Institute of Geoanthropology, Kahlaische Strasse 10, 07745 Jena, Germany

with highly structured content and complex layouts (Visually Rich Documents, VRDs) [1]. A high-profile example of this is the Habsburg Monarchy's *Hof- und Staatsschematismus*, which was published from 1702 to 1918 [2]. It provides us with a set of serial data on the administrative and representative elites of Habsburg Central Europe in very high quality [3]. A thorough analysis of this data set would contribute decisively to a significantly better understanding of relevant social processes, power dynamics, social networks and careers in modern Central Europe. The *Schematismus* could be used to trace the genesis of state and administrative institutions, their functioning and development, and the professional biographies of tens of thousands of officials and decision-makers over more than two centuries, across political ruptures and social transformations.

However, the complex structure of such publications has made a more comprehensive and quantitative evaluation of this source impossible. Even though OCR-quality has improved dramatically since the early days of the retrodigitisation of historical publications, structure analysis and layout detection remain a challenge. The *Schematismus* is known for its multi-column layouts, deeply branched hier-

Ministerium des kaiserlichen und königlichen Hauses und des Äußern. 289

<p>Lephanse Wenzel, <math>\Phi</math>C. päpstl. GO-R., Leut. i. d. R. Vize-Kons.</p> <p>Kobrasch Rudolf, <math>\Phi</math>C. Kons. Attaché.</p> <p>Wenzelauer zu Spermannsfeld Walter, v. v. Kons. Attaché.</p> <p>Wenzelka Avulin, <math>\Phi</math>C. Kons. Kanzlei-Schr.</p> <p>Untergeordnetes Amt: Konsulat.</p> <p>* In Durazzo.</p> <p>Halla Karl, FZO-R., <math>\Phi</math>C. pr. SLO. 2. Leut. im n. a. Stände der Landw., Vize-Kons. Gerent des Konsulates.</p> <p>Rutsky v. Rotok o. Divkafala Ludwig, <math>\Phi</math>C. J. Dr., Leut. i. d. Res., Vize-Kons.</p> <p>* Konsulat in Mitrovica.</p> <p>Taty v. Tatyev o. Tarkas Ludislav, <math>\Phi</math>C. Vize-Kons., Gerent des Konsulates.</p> <p>* Konsulat in Prizren.</p> <p>Prechaska Oskar, FZO-R., Leut. im n. a. Stände der Landw., Vize-Kons., Gerent des Konsulates.</p> <p>* Konsulat in Üsküb.</p> <p>Jurgatski Nikolai, RH. v. FZO-R., <math>\Phi</math>C. <math>\Phi</math>C. mem. DO 3. pr. HAO-R. 4., Leut. Leut.-Kons.</p> <p>Adamskiewicz Georg, Vize-Kons.</p> <p>Ritovsky v. Szepesszova u. Szoboraz Heinrich, <math>\Phi</math>C. Vize-Kons.</p> <p>Rozel v. Vukopas Tiber, Kons. Attaché.</p> <p>Dobratski Andreja, GVK. m. K., <math>\Phi</math>C. <math>\Phi</math>C. <math>\Phi</math>C. Kons. Kanzlei-Schr.</p> <p>* Konsulat in Monastir.</p> <p>Boronizza Julius, Freih. v. EKO-R. 3., <math>\Phi</math>C. Käm., Honvéd-Leut. i. d. Res., Vize-Kons.</p> <p>Zitovezi v. Szepesszova u. Szoboraz Heinrich, <math>\Phi</math>C. Vize-Kons.</p> <p>* General-Konsulat in Janina.</p> <p>Bilski Konstantin, <math>\Phi</math>C. <math>\Phi</math>C. ott. MO 3. ott. MO 4. Leut. i. d. Landw. Evidenz-Kons. m. d. I. a. Leg. Sekr.</p> <p>Molnar v. Mischkovicz Julius, <math>\Phi</math>C. ott. MO 4. Major d. R. Hon. Vize-Kons.</p> <p>Untergeordnete Ämter: Konsulat.</p> <p>* In Valona (Atona).</p> <p>Krusa Friedrich, <math>\Phi</math>C. Leut. i. d. Res., Vize-Kons., Gerent des Konsulates.</p> <p>Kouda Hugo, <math>\Phi</math>C. Kons. Kanzlei-Schr.</p> <p>Vize-Konsulat.</p> <p>In Prezera.</p> <p>Zaccaria A., <math>\Phi</math>C. Gerent des Vize-Konsulates.</p> <p>* General-Konsulat in Salonich.</p> <p>Pala Gottlieb, FZO-Kl., EKO-R. 3., <math>\Phi</math>C. <math>\Phi</math>C. päpstl. GO-R. m. St. H. KO-R. gr. EO-Kl. pr. KO-R. 2. r. AO-R. 2. ott. MO 2. boig. ZVO 3. Leut. i. d. Landw. Evidenz-Kons. m. d. I. a. Leg. Sekr.</p> <p>Gregorich Miklav, <math>\Phi</math>C. n. a. Landw. Leut., Vize-Kons.</p> <p>Filinger Hans, <math>\Phi</math>C. Leut. i. d. Res., Vize-Kons.</p> <p>Klocek Josef, GVK. m. K., <math>\Phi</math>C. <math>\Phi</math>C. <math>\Phi</math>C. Kons. Kanzlei-Schr.</p> <p>Untergeordnete Ämter: Vize-Konsulat.</p> <p>In Serres.</p> <p>Zlatko Georg C. FZO-R., <math>\Phi</math>C. gr. EO-R., boig. ZVO 4., Hon. Vize-Kons.</p>	<p>Konsular-Agentie.</p> <p>In Cavalla.</p> <p>Wie v. Zantva Adolf, <math>\Phi</math>C. pr. KO-R. 4., Hon. Vize-Kons. (ad pers.).</p> <p>* Konsulat in Adrianopol.</p> <p>Jezensky v. Kis-Jezens u. zu Felkafala Ludwig, EKO-R. 3., <math>\Phi</math>C. ott. MO 3. J. Dr., Kons. Kanzlei-Schr.</p> <p>Netrovich Edl. v. Castel-Tromb Matko, <math>\Phi</math>C. <math>\Phi</math>C. Kons. Kanzlei-Schr.</p> <p>Untergeordnete Ämter: Konsular-Agentie.</p> <p>Bergajlan Eugen, <math>\Phi</math>C. Hon. Kons. Agent.</p> <p>In Gallipoli (Osman, Helich).</p> <p>Siderides Theodoros, <math>\Phi</math>C. <math>\Phi</math>C. prov. Gerent der Kons. Agentie.</p> <p>In Kiraklisse.</p> <p>Dodepato D. Konstantin, <math>\Phi</math>C. Hon. Kons. Agent.</p> <p>In Porto Lagos (Xanthi).</p> <p>Bergajlan Eugen, Hon. Kons. Agent in Dodekanes, zugl. prov. Gerent der Kons. Agentie.</p> <p>In Rodosto.</p> <p>Astar Pierre, GVK. m. K., <math>\Phi</math>C. <math>\Phi</math>C. Hon. Kons. Agent.</p> <p>* Konsulat in Constantinopel.</p> <p>Pavliu Duido, FZO-R., <math>\Phi</math>C. <math>\Phi</math>C. ott. MO 4. Leut. i. d. Hon. Kons. Kanzlei-Schr.</p> <p>Franciahi Rudolf, v. <math>\Phi</math>C. <math>\Phi</math>C. Leut. i. d. Landw. Evidenz-Kons. m. d. I. a. Leg. Sekr.</p> <p>Hlavac Edl. v. Rechtwall Friedrich, <math>\Phi</math>C. Kons. Attaché.</p> <p>Invanich Silvio M. FZO-R., <math>\Phi</math>C. gr. VVO-R. Kons. Kanzlei-Dir., Hohen-Kaplan.</p> <p>Trand Ludwig, <math>\Phi</math>C. <math>\Phi</math>C. ott. MO 4., zugl. Hon. Drago-Kanzleiman. Sekr.</p> <p>Citterich Humbert, <math>\Phi</math>C. <math>\Phi</math>C.</p> <p>Untergeordnete Ämter: Vize-Konsulat.</p> <p>In den Bardanelien.</p> <p>Kantopulos Konstantin, FZO-R., <math>\Phi</math>C. <math>\Phi</math>C. FZO-R. 1. eidoch. HVO-Edl., ott. MO 3. Hon. Kons. (ad pers.).</p> <p>In Djedra.</p> <p>Tosid Dusan, <math>\Phi</math>C. M. Dr., Hon. Vize-Kons.</p> <p>Konsular-Agentie.</p> <p>* In Brussa.</p> <p>Stevens John, FZO-R., <math>\Phi</math>C. gr. EO-Off., ott. MO 4., Kons. Kanzlei-R., Gerent der Kons. Agentie.</p> <p>In Tenedos.</p> <p>Geraaglia Anton, GVK. m. K., <math>\Phi</math>C. <math>\Phi</math>C. Hon. Kons. Agent.</p>	<p>* General-Konsulat in Smyrna.</p> <p>Kral August, FZO-Kl., EKO-R. 3., <math>\Phi</math>C. <math>\Phi</math>C. päpstl. SO-Kl. m. St., pr. SLO 2. ott. MO 3. Kons. m. d. I. a. Gen. Kons. H. Kl., Gerent des Gen. Konsulats.</p> <p>Herzfeld Max, RH. v. <math>\Phi</math>C. <math>\Phi</math>C. ott. MO 3. J. Dr., Vize-Kons.</p> <p>Frossard Maxzeil Edl. v. Kons. Attaché.</p> <p>Vintara Remona, GVK. m. K., <math>\Phi</math>C. <math>\Phi</math>C. ott. MO 3. Kons. Kanzlei-Schr.</p> <p>Ferhat Viktor, GVK. m. K., <math>\Phi</math>C. <math>\Phi</math>C. ott. MO 4., Kons. Kanzlei-Schr.</p> <p>Untergeordnete Ämter: Vize-Konsulate.</p> <p>In Chios.</p> <p>Brazzafall Francesco, <math>\Phi</math>C. Hon. Vize-Kons.</p> <p>* In Rhodos.</p> <p>Barmann Anton, <math>\Phi</math>C. prov. Gerent des Vize-Konsulates.</p> <p>In Samos.</p> <p>Misar Oskar, <math>\Phi</math>C. ott. MO 4., Hon. Vize-Kons.</p> <p>Konsular-Agentie.</p> <p>In Metelin.</p> <p>Bargill Natalie, GVK. m. K., <math>\Phi</math>C. <math>\Phi</math>C. ott. MO 4., Hon. Kons. Agent.</p> <p>* Konsulat in Canoa.</p> <p>Wein Jakob, <math>\Phi</math>C. <math>\Phi</math>C. Honvéd-Leut. m. d. Res.</p> <p>Herzfeld Emserich, RH. v. <math>\Phi</math>C. Vize-Kons.</p> <p>Untergeordnete Ämter: Konsular-Agentie.</p> <p>In Candia.</p> <p>Teresio Viktor, Hon. Vize-Kons. (ad pers.).</p> <p>In Nettimo.</p> <p>Treffi Theodor, GVK. m. K., <math>\Phi</math>C. <math>\Phi</math>C. H. KO-R., gr. EO-R., ott. MO 4., Hon. Vize-Kons. (ad pers.).</p> <p>* General-Konsulat in Traperunt.</p> <p>Möriz v. Tescs Peter, <math>\Phi</math>C. <math>\Phi</math>C. Käm., Gen. Kons. H. Kl., Hon. ung. Kons. Oberrichter-Stellv.</p> <p>Untergeordnete Ämter: Vize-Konsulat.</p> <p>In Samos.</p> <p>Torre A. del, <math>\Phi</math>C. prov. Gerent des Vize-Konsulates.</p> <p>Konsular-Agentie.</p> <p>In Kerassunt.</p> <p>Alparzi G., <math>\Phi</math>C. prov. Gerent der Konsular-Agentie.</p> <p>* Konsulat in Aleppo.</p> <p>Pueho Friedrich, prov. Gerent des Konsulates.</p> <p>Untergeordnete Ämter: Konsular-Agentie.</p> <p>In Alessandretta.</p> <p>Levante Emil, <math>\Phi</math>C. <math>\Phi</math>C. Hon. Kons. Agent, m. d. I. a. Hon. Vize-Kons.</p> <p>In Herosia.</p> <p>Daras Nikolaus, FZO-R., <math>\Phi</math>C. gr. EO-R., Hon. Kons. (ad pers.).</p>
--	---	---

Fig. 1 Example page of *Schematismus* from 1910, highlighting the complexity of its structure, including multiple columns, hierarchical relationships, unique characters, and special annotations, such as the curly braces of which two of them are found in the middle of the page [4]

archies of several levels, and multimodal page designs that feature text alongside tables and complex lists. These characteristics represent the significant challenges digital historians and humanists encounter when processing large quantities of such documents. Thus, there has not been a comprehensive extraction of information, entire data sets or structures, such as information on hierarchies, on the detailed composition of administrative authorities, or simply careers or biographies with regard to the *Schematismus* yet. Some studies manually extracted relevant information, which proved tedious, time-consuming and prone to errors.

For years consideration has been given to publishing parts of the *Schematismus*-series as digital editions. Such undertakings have so far reliably failed because of the immense size of the task; the most important series of such handbooks, digitally published by the Austrian National Library<sup>1</sup>, comprises 145 volumes compiled between 1702 and 1918 (Table 1).

<sup>1</sup> alex.onb.ac.at.

Table 1 Data sheet of the *Schematismus*

Key statistics
Total of 145 volumes
56 volumes in the 1700s
89 volumes between 1800 and 1918
Complete series from 1816 to 1848 and from 1876 to 1918
Ranging from ~150 pages (1702) to ~1900 pages (1918)

We calculated that this series contains between 130,000 and 150,000 printed pages. Neither manual extraction of the information contained therein nor automated processing of the pages, which have very diverse layouts, seemed feasible to us with the solutions currently available. The off-the-shelf solutions we tried, as well as generic layout detection integrated into well-established OCR did not perform too well with the *Schematismus*.

We assumed that a powerful layout detection, which can divide the individual cells of the *Schematismus* into meaningful text blocks, could possibly represent a relatively favourable solution. The complexity and diversity of the layouts led us to consider machine learning models as possible solutions, as we expected them to be able to handle the fractal nature of the page layout better than rule-based models. The first research question in this paper therefore investigates, whether high-quality layout detection as a preprocessing step could improve the performance of downstream OCR, in order to obtain a relatively simple solution for extracting the relevant information. We first identified and tested a suitable deep learning architecture—Faster R-CNN. Then we rebuilt the custom font used in the *Schematismus*. We used this font to synthesize a large amount of *Schematismus*-styled training data in the first step. Care was taken to ensure that the synthetic training data had a similarly complex and varied layout as the original data sets. The synthetically produced training data was also artificially distorted, dirtied, and twisted. We trained the Faster R-CNN model with the synthetic training data and further finetuned the model with a smaller number of manually annotated pages from the *Schematismus*.

Once our layout detection was working, we tested its potential with an off-the-shelf distribution of Tesseract.

Then we addressed the second research question, which investigated, to what degree a custom font finetuned distribution of Tesseract could boost OCR results. We finetuned an off-the-shelf distribution of Tesseract with the custom font we had created of the *Schematismus* and carried out a second test run.

For this study, we only invested limited resources in comparing possible OCR solutions, and we only took a cursory look at the necessary post-processing steps. However, the layout detection solution we present should be able to operate with any downstream OCR solution that can be integrated

into a Python pipeline. We expect that it might perform even stronger with OCR-engines developed specifically for historical fonts, such as Kraken OCR or OCR4all.

We have not devoted significant resources to the post-processing step so far, yet we ran some preliminary tests, using GPT4. The results were promising, and we expect that the latest generation of LLMs provides enormous potential to further improve the quality of the OCR results in a post-processing step.

## 2 Background & related work

Wide variation in document formats and degradation over time make automated layout detection and the development of generic tools quite challenging, especially for historical documents. For these reasons, a broad range of techniques and methods has been developed and applied over time, in order to be able to provide for a better automated information extraction from historical documents.

In this section, we will first present the state of research in history and more broadly the humanities, then initiatives to improve OCR that have been developed and tested in the field of historical document analysis are discussed.

### 2.1 Availability and limitations of OCRed primary sources

The retro-digitisation of either parts of or entire historical sources has been an issue among historians and fellow humanists since at least the 1950s, particularly with regard to serial sources [5]. However, for the larger part of the past seven decades, information extraction and the production of digital data that could be processed by computing machines, has been executed manually. Even though OCR software has been widely available by the 1990s at latest, particularly processing historical data has remained a complex issue [6], as the quality of results has been varying strongly [7]. We therefore observe a bifurcation in the field of digital historical and humanist research: On the one hand, large amounts of retro-digitised historical data are processed automatically by important providers of research data, such as [www.archive.org](http://www.archive.org) [8]. Many national libraries, as for instance the Austrian National Library [9], the National Library of Finland [10] or the Munich Digitization Center of the Bavarian State Library [11], and further large transnational initiatives, as for example Europeana [12], also belong to this group. However, users of these data, such as historical research projects, are still relying on the manual transcription or annotation of digitised primary sources, even at scale.<sup>2</sup> Frequently, this is

<sup>2</sup> There is very little documentation on how data in historical research projects is OCRed, so we have to refer to information that we obtained in

due to quality issues regarding OCRed documents available online.<sup>3</sup> Relevant databases that were build predominantly on data extracted by manual labour include the *Wiener Datenbank zur Europäischen Familiengeschichte* [13], *The Emperor's Desk* [14], the prosopographical data processed by *The Viennese Court* [15], and further the projects run in *Social Mobility of Elites* [16].

### 2.2 Common challenges

Historians and digital humanists working with retro-digitised text data are familiar with the phenomenon that OCR-ed texts provided on common platforms may vary in quality. This is mostly due to the fact that texts underwent OCR procedure at different times and different technologies were used. As a result, full-text searches frequently produce inhomogeneous results. Also, a “systematic” OCR error, i.e. distortions that are typical for certain OCR software, can no longer be processed in such a targeted manner if large text corpora consist of parts that were OCRed at different times with different software. In many projects, therefore, raw digitised documents are now re-OCRed, although many historical layouts still pose a challenge for standard OCR but also specialised software [10].<sup>4</sup> Apart from specialised solutions [17], Tesseract OCR and Transkribus are currently regarded as reference standards in the field of historical OCR, but still encounter limitations that require a high level of manual effort, given the special layouts and structures that will be discussed here [18]. Further, OCR4all offers a toolbox for open source OCR applications that has proved highly performant lately [19], as well as kraken OCR, developed by the EHESS and closely associated with the digital research environment eScriptorium.<sup>5</sup> Libraries have begun to use the potential of ML to the classification of large quantities of texts [20], especially in the context of research libraries [21].

For historical research, important primary sources such as censuses of the Habsburg Empire have been mostly digitized manually [22, 23]. Such sources usually feature a relatively limited volume and a highly homogeneous layout and struc-

personal conversation with several dozen colleagues over the past one-and-a-half decade. Based on the information available to us, we would assume that manual transcription of historical text, but especially tabular data, is still the standard procedure, even though this usually starts with a raw document created with standard OCR (Tesseract or Transkribus), which is subsequently improved. The digitisation of tabular data is particularly challenging because the target structure usually has to be created manually beforehand. Cf. [6].

<sup>3</sup> An entire panel in the recent 15<sup>th</sup> Austrian Contemporary History Day was devoted to the manual and semi-automatized correction of flawed OCR data from large repositories.

<sup>4</sup> Details about the OCR process at large providers of cultural heritage data are only available to a very limited extent, and the publicly accessible documentation leaves much to be desired, even at public institutions.

<sup>5</sup> <https://github.com/mittagessen/kraken>.

ture. However, because of the enormous volume, this method is not feasible for other complex source works, such as the *Schematismus*. Even in the last initiative we know of, the size and complexity of the *Schematismus* was ultimately considered insurmountable for only partially automated data extraction, and manual edition of a small part was envisaged as an alternative. Due to these obstacles, very little research is engaging with a deeper exploration of the information stored in the *Schematismus*, the work of Bavouzet [24] (building entirely on manually extracted data) is clearly standing out as a beacon here.

Only recently have developments in new fields of research such as document intelligence begun to open up the possibilities of machine learning for this area on a larger scale [25–27]. Impressive progress has been made in some areas [28].

### 2.3 Pre-OCR steps to enhance OCR quality

The process of extracting information from historical documents can also be considered as a complete extraction pipeline, instead of individual tasks. Monnier and Aubry present work on an extraction tool for historical documents, which provides improvements in terms of robustness and extraction performance due to mutual reinforcement of text line and image segmentation [29]. The task of segmentation is also highlighted by Gruber et al. [30] where the authors also propose to conduct preprocessing of the image before conducting OCR. Such approaches have already been explored in the past, for example by processing the background of the image [31]. Consequently, techniques such as Generative Adversarial Networks have been explored to achieve super-resolving of the input images [32]. Augmentation has been used on several occasions lately in order to develop economic ways to scale up information extraction from historical documents. Grüning et al. use the deep neural network ARU-Net to address the issue of line detection in historical documents [33]. Martínek et al. [18] realise an approach, which combines fine-tuning OCR-engines with comparatively little data, after training the engine with large synthetic data-sets.

Document layout analysis has received increasing attention in the past few years, though this area remains under-explored [34]. Solutions found in this field are not always tackling specifically historical problems, as for instance [35]. It has been recognised though that the increasing availability and usability of deep neural networks, in particular CNNs, offers entirely new opportunities for the development of custom-made solutions for certain document layout analysis tasks [36]. With regard to OCR, line detection is frequently considered more relevant than layout detection [28]. Nevertheless, significant progress has been made in this area in recent years, with layout detection often being understood as part of a more complex single stage process [1, 26, 27]. LayoutLM presents a promising and versatile solution to address

**Table 2** Approximate number of persons mentioned in each *Schematismus* per year and the number of pages in the respective volume

Year	Persons	Pages
1702	2200	150
1750	4000	430
1801	6700	1030
1820	26,400	1300
1848	28,800	1200
1860	15,000	300
1878	51,200	1100
1910	93,400	1600
1918	110,400	1900

We calculated these values by averaging the number of bounding boxes predicted on each name register page and multiplying this average by the number of name index pages listed in each year's name register

many common tasks in this area [37]. We also tried LayoutLM for the *Schematismus*, but it did not prove efficient with this particular type of documents.

## 3 Methodology

State Manuals and *Hof- und Staats-schematismen* for the Habsburg Empire in particular are commonly provided in PDF format. These documents are accessible via various historical document repositories such as the Austrian National Library<sup>6</sup>, the Munich Digitalization Center (MDZ) of the Bavarian State Library<sup>7</sup> or archive.org<sup>8</sup>. Frequently, however, plain text is not available at all or is available in sub-optimal quality.

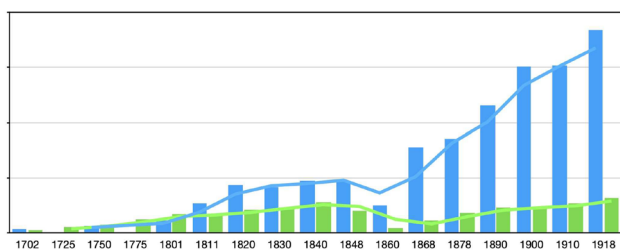
Whereas the general issue of suboptimal OCR quality is widely acknowledged in historical and further humanist research, most approaches consider layout a secondary line of attack, when it comes to improving OCR quality. There are few exceptions, as for instance [34, 36, 38]. Even the two probably most important and most widely used out-of-the-box solutions in the field of historical OCR, different distributions of Tesseract and the variety of different OCR and HTR models provided by Transkribus, require significant (manual) effort in data preprocessing, when it comes to extracting information from digitised primary sources. Even conventional document layouts often pose a challenge [6]. Lately, OCR4all has come up with a strong layout detection.

Complex layouts as encountered in the *Schematismus*, however, are still considered a major challenge by the his-

<sup>6</sup> [alex.onb.ac.at](http://alex.onb.ac.at).

<sup>7</sup> [www.digitale-sammlungen.de](http://www.digitale-sammlungen.de).

<sup>8</sup> [www.archive.org](http://www.archive.org).



**Fig. 2** Number of persons (blue) and pages (green, 10fold augmentation) in different volumes of the *Schematismus*

torical research community. Efficient information extraction from such documents is considered a difficult and complicated task. Several efforts to automatically extract structured data from the *Schematismus* have failed so far. Generally, particularly in historical texts, OCR is performed page-wise, which means that snippets, belonging to different blocks of text, are performed serially, as common OCR processors work line-wise.

### 3.1 Approach and research questions

For our approach, we identified two points of attack, which correspond with the two research questions we formulated. First, we wanted to find out, whether AI-driven layout detection prior to OCR could significantly boost OCR quality. Second, we were interested to see, to what degree finetuning the OCR engine could further improve its performance. To tackle research question 1 in a first step, we split the individual document pages into their layout elements, in order to preserve the context of the different blocks of text. To this end, we used a deep learning convolutional neural network. This is a complex ML algorithm originally developed for object localization and object recognition. We assumed this approach would be suitable, since we consider the identification of large coherent blocks a computer vision problem.

In the next step, we used an OCR algorithm to process the individual image snippets, rather than the entire document page image. Then we addressed the second research question, to which degree OCR accuracy could be improved by finetuning the standard OCR tool Tesseract on a custom font, which was designed to look as similar as possible to the original font used in the *Schematismus* documents.

We chose a sample from 145 volumes of State Manuals and focused on editions that were published in the second half of the 19th century onward, for two reasons:

1. The task is becoming slightly more complex for the decades prior to 1848, as the fonts that were used are more diverse and complex. We do not consider this a major problem, yet for the proof of concept we were interested in streamlining the entire research process and

to eliminate additional complexities that were not in our primary line of attack.

2. Even though State Manuals are available for a period of more than 200 years, the mass of data was produced from the 1850s onward, therefore the yield for a solution capable of dealing with documents of this type is expected to be very high (compare Fig. 2 and Table 2).

To deal with research question 1, we built a machine learning model that can segment retrodigitised PDF-documents of the *Hof- und Staatshandbücher* and split these into their layout-structure elements, such as individual paragraphs and headings. Each of these image snippets was subsequently fed into Tesseract for text extraction. Figure 3 shows a simplified version of this process.

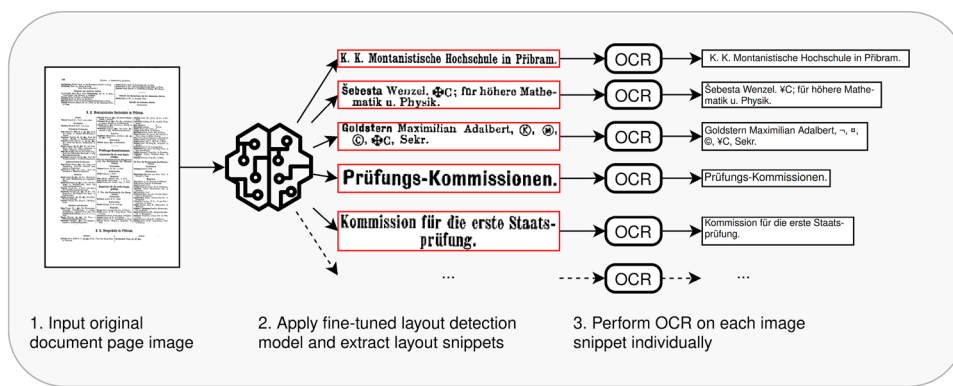
We used two Tesseract OCR models. For research question 2, an instance of Tesseract that had been fine-tuned on our custom font. For comparison and performance assessment, we also used an instance that had not been fine-tuned. OCR accuracy was then calculated by comparing the extracted text with the manually transcribed ground truth. Then, this process was repeated, but without dividing the page into individual segments. To answer the research questions in this study, the accuracies were finally compared in order to evaluate the efficiency of our approach.

The following subsections will describe in detail the implementation of the different methods that we employed to process PDFs of the *Hof- und Staatshandbücher*. This includes training data set generation for the development of the layout detection, layout detection itself, and optical character recognition. Each of these three steps constituted a work package of its own. All the research and analysis was conducted within Jupyter Notebook, the diverse tools that were put to use are listed in the following subsections.

### 3.2 Data set generation

In order to successfully train a convolutional neural network, a sufficient amount of labeled data is required. Creating a training data set by manually drawing bounding boxes on a large quantity of PDFs drawn from historical source documents is time-consuming and labour-intensive due to the large number of pages that would have to be annotated. Therefore, we developed an alternative approach to artificially generate labeled training data. We wrote a Python script that is designed to generate synthetic documents mimicking the style of the *Hof- und Staatshandbücher*. This Python script generated Latex-code, which was then compiled using `luatex` [39] to create a PDF file. In the course of this process, the coordinates of the individual text structure elements were determined, which is described in more detail in the following section.

**Fig. 3** This flowchart illustrates how each extracted layout element is processed by OCR



In order for the generated data to be used as training data, it was imperative that the created documents appeared as visually realistic as possible compared to the original documents. Reverse engineering the original documents and paying attention to detail were therefore essential to the creation of synthetic training data.

Due to the fact that each paragraph begins with a last name and a first name, a data set containing thousands of first and last names [40] was randomly sampled in order to obtain names. The pool of samples was restricted to Austria, Hungary, Switzerland, and Germany.

A list of abbreviation explanations from the 1910 *Schematismus* document was manually transcribed and randomly selected for the text following the names. To vary the length of individual generated paragraphs, the number of sampled texts was also randomly chosen.

Figure 1 displays a page from the original 1910 *Schematismus* document. Variable column numbers are a key characteristic of such documents and are used almost everywhere. While most sections have three columns, there can be variations in certain sections. For instance, Fig. 8 shows four columns in the name-index section. Thus, generating realistic synthetic documents required the use of the same column layout.

Another key visual element is the relatively distinct font type. Research led us to a font called “Opera-Lyrics-Smooth” that appeared very similar to the original. Even though it represented already a good match, we decided to invest additional effort: Using the open source program “FontForge” [41], we further customised the font in accordance with the original. In order to achieve the best possible match to the original font, screenshots of every letter in the original documents were taken manually, and then the existing letters in the font were adapted according to the screenshots taken.

Another distinctive feature of the *Hof- und Staats-schematismus* is its excessive use of particular symbols, representing orders and similar distinctions of the persons listed in this source. The use of these symbols allowed the further condensation of information stored in the *Schematismus*.

Unicode	Symbol
¢	Ⓢ
§	Ⓢ
£	Ⓢ
»	Ⓢ
¥	Ⓢ
!	Ⓢ
-	Ⓢ
⊗	Ⓢ
©	Ⓢ
«	Ⓢ
┌	Ⓢ

**Fig. 4** Illustration of how Unicode characters are mapped to *Schematismus* symbols

Nonstandard symbols were mapped to Unicode symbols, which were unlikely to be needed for document generation. Figure 4 illustrates this mapping.

This method resulted in the creation of three font types: one for general text paragraphs, one for headlines, and one for italics. An example of this can be seen in Figs. 5 and 6, which shows an original paragraph and a corresponding paragraph reproduced using the custom font. Additionally, when reviewing some original *Schematismus* documents, it is apparent that font sizes and alignments vary considerably from section to section or even from page to page. In particular, the difference can be observed when examining the headlines of the original documents. To create realistic head-

**Wetschl Franz, Freih. v., FJO-GK., EKO-R. 2., LO-R.,** ④, ①, ②, ③, ④, ⑤, ⑥, ⑦, ⑧, ⑨, ⑩, ⑪, ⑫, ⑬, ⑭, ⑮, ⑯, ⑰, ⑱, ⑲, ⑳, ㉑, ㉒, ㉓, ㉔, ㉕, ㉖, ㉗, ㉘, ㉙, ㉚, ㉛, ㉜, ㉝, ㉞, ㉟, ㊀, ㊁, ㊂, ㊃, ㊄, ㊅, ㊆, ㊇, ㊈, ㊉, ㊊, ㊋, ㊌, ㊍, ㊎, ㊏, ㊐, ㊑, ㊒, ㊓, ㊔, ㊕, ㊖, ㊗, ㊘, ㊙, ㊚, ㊛, ㊜, ㊝, ㊞, ㊟, ㊠, ㊡, ㊢, ㊣, ㊤, ㊥, ㊦, ㊧, ㊨, ㊩, ㊪, ㊫, ㊬, ㊭, ㊮, ㊯, ㊰, ㊱, ㊲, ㊳, ㊴, ㊵, ㊶, ㊷, ㊸, ㊹, ㊺, ㊻, ㊼, ㊽, ㊾, ㊿, tsc. ZVO-GOff., r. StO-R. 1., prs. SLO. 1., sp. IO-GK., rm. KO-GK., pr. RAO-R. 2. m. St. i. Br., pr. KO-R. 2. m. St., bayr. KO-GKt., s. AO-Kt. 1., wt. KO-Kt. m. St., schw. NSO-Kd. 1., blg. LO-GOff., siam. KO-GOff., SEHO-Kt. 1., Greffier des österr. kais. Leopold-Ordens, Mitgl. der Stadterweiterungs-Kmsn. im Mstm. des Innern.

**Fig. 5** Original paragraph, taken from an existing state manual from the year 1910

**Wetschl Franz, Freih. v., FJO-GK., EKO-R. 2., LO-R.,** ④, ①, ②, ③, ④, ⑤, ⑥, ⑦, ⑧, ⑨, ⑩, ⑪, ⑫, ⑬, ⑭, ⑮, ⑯, ⑰, ⑱, ⑲, ⑳, ㉑, ㉒, ㉓, ㉔, ㉕, ㉖, ㉗, ㉘, ㉙, ㉚, ㉛, ㉜, ㉝, ㉞, ㉟, ㊀, ㊁, ㊂, ㊃, ㊄, ㊅, ㊆, ㊇, ㊈, ㊉, ㊊, ㊋, ㊌, ㊍, ㊎, ㊏, ㊐, ㊑, ㊒, ㊓, ㊔, ㊕, ㊖, ㊗, ㊘, ㊙, ㊚, ㊛, ㊜, ㊝, ㊞, ㊟, ㊠, ㊡, ㊢, ㊣, ㊤, ㊥, ㊦, ㊧, ㊨, ㊩, ㊪, ㊫, ㊬, ㊭, ㊮, ㊯, ㊰, ㊱, ㊲, ㊳, ㊴, ㊵, ㊶, ㊷, ㊸, ㊹, ㊺, ㊻, ㊼, ㊽, ㊾, ㊿, tsc. ZVO-GOff., r. StO-R. 1., prs. SLO. 1., sp. IO-GK., pr. RAO-R. 2. m. St. i. Br., pr. KO-R. 2. m. St., bayr. KO-GKt., s. AO-Kt. 1., wt. KO-Kt. m. St., schw. NSO-Kd. 1., blg. LO-Goff., siam. KO-Goff., SEHO-Kt. 1., Greffier des österr. kais. Leopold-Ordens, Mitgl. der Stadterweiterungs-Kmsn. im Mstm. des Innern.

**Fig. 6** Corresponding synthetically generated paragraph reproduced with a custom font

lines, four headline types ranging from “H1” to “H4” were used to emulate this feature. The first element was the largest and the rest gradually decreased in size. Table 3 provides a list of all the different classes.

Finally, it is crucial to emphasise some small but very significant visual details. Every paragraph begins with one or two words in bold, the last name of the individual, followed by the first name and some additional titles and awards. Indentation occurs after the first line if the text is too long for a single line. Furthermore, multiple individuals may be grouped together within a large curly bracket, as can be seen in Fig. 1. In name-index pages, multiple individuals with the same last name may be grouped by adding a horizontal line at the beginning, which can be observed in Fig. 7. Additionally, every entry within this section is accompanied by one or more numbers that indicate the page number. Finally, it should be noted that headings are usually centered on the page or within columns and that the end of every text is always marked with a period.

In order to be able to produce text for the synthetic *Schematismus*-style training documents, several sources were consulted. For the purpose of generating large headlines, a simple list of historical Austrian orders and decorations was used [42]. In order to create headlines with smaller font sizes, a combination of years as strings and Austrian municipality names was used. For paragraph generation, two sources were consulted, as previously mentioned. Through the use of all the above methods and visual keys, we were able to create a large number of realistic looking synthetic docu-

ments. An example of such a synthetic *Schematismus*-style document can be seen in Figure 7.

In addition to generating a synthetic data set of *Schematismus* documents for training a machine learning model, annotations for each element of the structure had to be generated along with the generation of the document in order to make this data set effective for training. As part of the process of detecting and localising objects, in this case the layout elements, bounding boxes were used to define the location and size of the individual structure elements within an image. The labels accompanying the bounding boxes indicate the class of the corresponding box, such as “paragraph” or “H1”, which are necessary for classification tasks performed by the machine learning method.

A latex package called “zref-savepos” is used to save the position of characters on the current page and write these coordinates to an external file at compilation time. Using the coordinates, bounding boxes could be computed by parsing the external file. As the individual text elements had already been generated earlier in the same Python script, it was known which label had to be associated with the respective bounding box. In order to construct the data set, the generated documents, which were compiled by the latex compiler and then saved as PDF files, were converted to images. In addition to the image file, a Pascal VOC XML file containing the corresponding annotations was created and stored in a separate directory. A total of 3766 synthetic *Schematismus* documents have been generated using this approach. Figure 9 shows a synthetic document with its corresponding annotations overlaid.

### 3.3 Layout detection

For the actual layout detection model to be configured in the next step, we chose a faster region-based convolutional neural network (faster R-CNN) built on a ResNet-50 backbone. The model was created and trained using the PyTorch [43] framework. In PyTorch, version-2 of the faster R-CNN implementation was used [44]. Training of the model was conducted on an Nvidia RTX 4090 with 24 GB of video memory.

#### 3.3.1 Model training and settings

Even though we used primarily the default settings of the model, we found that some adjustments had a significant impact on the model’s performance. The following paragraphs will provide a detailed description of these adjustments.

1. By setting the pre-trained parameter to “True”, the training speed has improved in a way that earlier epochs of training have already begun with a lower training and

**Fig. 7** Example of a synthetic Schematismus-style document used in training-set

<b>Marianerkreuz des Deutschen Ritterordeas</b>		
<p>Kriegskreuz für Zivilverdienste II. Klasse.</p> <p><b>Ratz</b> Martina, Inf., R., LO., 1864-M., tsc. ZVO.</p> <p><b>Lisa</b> Luisa, prov., DevK., StKO., 1864-M., d. DO.</p> <p><b>Weber</b> Dennis, ER., ElisO., Albr-Z., VO. d. bayr. Kr.</p> <p style="text-align: center;">Geistliches Verdienstkreuz II. Klasse.</p> <p><b>Mouthon</b> Francine, Kord., FJO., @, hz. HO.</p> <p style="text-align: center;">Ehrenritterkreuz des Souveränen Malteserordens.</p> <p><b>Hübner</b> Peter, GK., FJO., @</p> <p><b>Armillotta</b> Nunzia, Bez., Kt., EKO. 3., GTM.</p> <p><b>Fuchs</b> Eric, prakt., OffK., EKO. 1., D1</p> <p><b>Porte</b> De-id. techn., EK., StKO., @</p> <p><b>Meyer</b> Torsten, m. E., StphO., G<sup>TM</sup>., parm. LO.</p> <p><b>Nüchter</b> Olaf, m. Schl., StphO., D2.</p> <p><b>Thiele</b> Marco, M., EKO. 3., @</p> <p><b>Wagner</b> Stefan, i. Br., StphO., TLVM.</p> <p><b>Ness</b> Ness, Kgs., V., StphO., STM. 1. 2.</p> <p><b>Jens</b> Weine, Lehrer-Bild., Anst., m. E., ElisO., @, prt. MVO.</p> <p><b>Seiner</b> Lila, ak., DonK., ElisO., @</p> <p><b>Domenico</b> Quattrocchi, Kr., Ger., GBd., EKO. 1., SVK.</p> <p><b>Mair</b> Kathrin, Goff., DRO., @</p> <p><b>Knudsen</b> Friedrich, Mag., i. Br., StphO., @, it. KO.</p> <p><b>Güral</b> Sara Meliss, V., FJO., D2., prt. BAO.</p> <p><b>Noldes</b> Lukas, GKord., LO., @, VO. d. bayr. Kr.</p> <p><b>Schrath</b> Hubert, Rgt., R., ElisO., AhWN, D2., bayr. KO.</p> <p><b>Rei</b> Arne, EZ., EKO. 1., @C</p> <p><b>Finke</b> Robert, kön., Ekt., MThO., @</p> <p><b>Stross</b> Harald, DonK., StKO., @, mx. AO.</p> <p><b>Hartmann</b> Christa, Hus., Husaren., GBd., EKO. 3., GTM., rm. StO.</p>	<p style="text-align: center;">mit der Krone.</p> <p style="text-align: center;">Ritterkreuz des Militär Maria-Theresien-Ordens.</p> <p><b>Kidane</b> Asia, Rgs., Kr., DRO., MKDRO., prt. VVO.</p> <p><b>Herbst</b> Heiko, A. B., GKord., MThO., @</p> <p style="text-align: center;">Vom Jahre 1874.</p> <p><b>Niemeier</b> Felix, F. Z. M., EB., EKO. 3., 1864-M.</p> <p><b>Wegner</b> Volker, m. Schl., EKO. 2., AhWN.</p> <p><b>Kremer</b> Stefan, EK., DRO., MKDRO.</p> <p><b>Galitzdorfer</b> Sandro, m. St., LO., @</p> <p><b>Smolle</b> Michaela, Insp., DevK., MThO., @, tun. NIO.</p> <p><b>Kaya</b> Bedri, i. Br., MThO., GVK., bd. ZLO.</p> <p><b>Nagy</b> István, V., GVIO., gsl. VK.</p> <p><b>Reinhard</b> Deroo Diane, m. L., FJO., ElisM., lip. HO.</p> <p><b>Schindhelm</b> Dennis, GK., EKO. 1., AhWN., pr. O. p. 1. m.</p> <p><b>Savi</b> Goran, m. K., EKO. 2., GVK., venezol. LibO.</p> <p><b>Heller</b> Nico, m. E., EKO. 3., @, chin. DO.</p> <p style="text-align: center;">Insigne des in Tirol immatrikulierten Adels.</p> <p><b>M Aldin</b> Mohamed, DevK., MaltO., @</p> <p><b>Thierauf</b> Michael, Lehen-Allod. Zentr. Kmsn., GK., EKO. 3., MfWK.</p> <p><b>Jethro</b> Leroy, Kd., DRO., @</p> <p><b>West</b> Anja, m. Schl., FJO., 1864-M.</p> <p><b>Teichmann</b> Jonas, R., Gkt., ElisO., Albr-Z., bulg. FO.</p> <p><b>Caselowsky</b> Carola, Min., GKord., GVIO., @, fz. MLO.</p> <p><b>Schaubach</b> Gerhard, Hus., Husaren., M., EKO. 3., EZfKW., est. AO.</p> <p style="text-align: center;">Ehrenritterkreuz des Souveränen Malteserordens.</p> <p><b>Born</b> Thomas, Gouv., EB., FJO., STM. 1. 2., tsc. ZVO.</p> <p style="text-align: center;">Erinnerungszeichen für die Ritter vom Goldenen Sporn.</p> <p><b>Maria</b> Mario, m. Schl., ElisO., @</p> <p><b>Meuter</b> Jakob, Prof., Kmsn., gld., StphO., STM. 1. 2.</p> <p><b>Aydin</b> Betül, M., GVIO., 1864-M.</p> <p><b>Montanus</b> Marion, (KD.), GVIO., @</p> <p><b>Glittenberg</b> Stefan, Assist., GK., EKO. 1., GTM., mekl. KO.</p> <p>□□□□ □□□□ □□□□ a. ö., m. K., StphO., MK-DRO., gr. EO.</p> <p style="text-align: center;">Ehrenzeichen für Verdienste um das Rote Kreuz I. Klasse.</p> <p><b>Szozurtek</b> Halina, m. K., GVIO., @</p> <p style="text-align: center;">Eisernes Verdienstkreuz.</p> <p><b>Lausanne</b> Vikings, math., GKord., StKO., D1, mekl. GO.</p> <p style="text-align: center;">Vom Jahre 1909.</p> <p style="text-align: center;">Jubiläums-Hof-Medaille.</p> <p style="text-align: center;">Militärdienstzeichen I. Klasse für Mannschaften.</p> <p><b>Steiner</b> Jeannette, Ekt., EKO. 3., STM. 1. 2.</p> <p><b>Kleinfeldt</b> Carolin, Übgssch. Lehrer., GK., MThO., ElisM.</p> <p><b>Acar</b> Halil, m. L., EKO. 2., @, siz. FdO.</p> <p><b>Lippens</b> Milena, gld., EKO. 1., MfWK., päpstl. EK.</p> <p><b>Marquardt</b> Ramona, EB., StKO., @, han. EAO.</p> <p><b>Latendin</b> Thomas, Kr., EKO. 3., GTM.</p> <p><b>Bartak</b> Andreas, GKord., FJO., D1, ix. MZVO.</p>	
<p style="text-align: center;">Bronzene Militär-Verdienstmedaille.</p> <p><b>Wisny</b> Halina, F. M. L., OffK., EKO. 1., MVK., jap. ChO.</p> <p><b>Nitschke</b> Daniela, m. Schl., EKO. 1., @, rm. StO.</p> <p><b>Bestmann</b> Andrea, Dion., Kt., FJO., @, pr. KO.</p> <p><b>Bernhard</b> Klaus, R., MThO., @H, VO. d. pr. Kr.</p> <p><b>Marelli</b> Stefano, GKord., EKO. 1., @, prt. VO.</p> <p style="text-align: center;">Kriegskreuz für Zivilverdienste IV. Klasse.</p>	<p style="text-align: center;">Kriegskreuz für Zivilverdienste II. Klasse.</p> <p><b>Ahren</b> Samuel, Konsist. R., M., StKO., ElisM., ndl. WO.</p> <p><b>Béres</b> Georgina, Magnatenh., Kr., EKO. 2., MfWK., bd. TrO.</p> <p><b>Reed</b> Dean, m. Schl., EKO. 1., Albr-Z.</p> <p><b>Kruber</b> Kerstin, Abg., Gkt., StphO., MfWK., s. RKO.</p> <p><b>Amort</b> Alexander, GKord., DRO., @</p> <p style="text-align: center;">Silbernes Verdienstkreuz</p>	<p><b>Koltai</b> Laszlo, gld., StKO., @, bayr. MxO.</p> <p style="text-align: center;">Ehrenzeichen für Verdienste um das Rote Kreuz II. Klasse.</p> <p><b>Yildiz</b> Yukseil Yusemin, G. M., OffK., MaltO., MKDRO., prt. MChO.</p> <p><b>Esai</b> Heidi, Gymn., ER., MThO., @, s. RKO.</p> <p><b>Westmark</b> Uwe, Stath., Ekt., DRO., @</p> <p><b>Loose</b> Stephanie, Suppl., EK., GVIO., @, brsch. HLO.</p> <p><b>D'sei</b> Celik Gonul, (KD.), StKO., D3.</p> <p style="text-align: center;">Vom Jahre 1904.</p> <p style="text-align: center;">Ehrenritterkreuz des Deutschen Ritterordens.</p> <p><b>Maser</b> Melitta, gld., StphO., @, bayr. MO.</p> <p><b>Kovács</b> Levente, A. K., Kd., EKO. 3., SVK.</p> <p><b>Marchlewski</b> Hans Peter, m. K., DRO., EZfKW., wt. OO.</p> <p><b>Naome</b> Ni, Kl., m. E., DRO., MfWK., schw. KO.</p> <p style="text-align: center;">Silberne Militär-Verdienstmedaille.</p> <p><b>Bakroné</b> Orsolya, n. a., Kt., GVIO., @C, hess. PhO.</p> <p><b>Mattern</b> Simone, wiss., EZ., MaltO., @</p> <p><b>Melzner</b> Sascha, polit., R., MaltO., Albr-Z., SMarO.</p> <p style="text-align: center;">Vom Jahre 1871.</p> <p><b>Ott</b> Wolfgang, GK., ElisO., GTM., tsc. JO.</p> <p><b>Zehnder</b> Damodar, R., EKO. 3., SVK.</p> <p><b>Bischof</b> Wolfgang, gld., EKO. 1., MVK.</p> <p><b>Bücker</b> Britta, R., FJO., @</p> <p><b>Njegova</b> Anđjelina, DonK., EKO. 3., 1864-M., prt. M<sup>VO</sup>.</p> <p><b>Lenz</b> Judith, GK., ElisO., GTM.</p> <p><b>Eich</b> Marion, OffK., EKO. 1., GVK.</p> <p><b>Huebner</b> Volker, R., FJO., ElisM.</p> <p><b>Ismael</b> Timiro, EK., MaltO., @</p> <p><b>Barleben</b> Rene, Goff., DRO., @</p> <p><b>Sportler</b> Florian, theol., m. E., MThO., TLVM., pr. KO.</p> <p><b>Sielaff</b> Andre, kan., m. St., EKO. 2., AhWN., han. GO.</p> <p style="text-align: center;">Großkreuz des Franz-Joseph-Ordens.</p> <p><b>Wilke</b> Majk, Geburtsh., m. E., StKO., @</p> <p><b>Jadon</b> Mostafa, m. Schl., LO., MKDRO., fz. HGO.</p> <p><b>Blastoek</b> Kirsten, Goff., DRO., MfWK.</p> <p><b>Deggar</b> Wiebke, Ekt., DRO., @, mekl. KO.</p> <p><b>Fritzen</b> Angelika, Volkssch., (KD.), GVIO., SVK., päpstl. PO.</p> <p><b>Kirkunki</b> Salar, Univ., Kd., EKO. 2., @</p> <p><b>Tak</b> Abi Tak, m. St., ElisO., ElisM., ix. EKO.</p> <p style="text-align: center;">Ehrenritterkreuz des Souveränen Malteserordens.</p> <p><b>Wagner</b> Arno, M., EKO. 1., @H</p> <p><b>Eberle</b> Hans-Jörg, Kd., ElisO., EZfKW.</p> <p><b>Arshgesicht</b> Patrik, Goff., StKO., @</p> <p><b>Wohletz</b> Matthias, Stadt R., Gbd., MaltO., ElisM.</p> <p><b>Pieck</b> Joel, Ther. Ak., EZ., StKO., @, päpstl. ChO.</p>

validation loss than with randomly initialised starting weights.

2. A further parameter that has been tweaked is the anchor generation of the region proposal network in the faster R-CNN model. When identifying areas of interest in an image, anchor boxes are critical because they determine where to look. Thus, they play an essential role when it comes to detecting layout elements within a document image. In order to be able to accommodate a variety of different types of objects, anchor boxes were selected with different aspect ratios and scales. It is imperative

to note that there are multiple anchor boxes applied to each sliding window position in the region proposal network. Therefore, it is logical to specify these ratios and scales according to the shape of the objects. Thus, the minimum, maximum, and mean ratio and scale of all bounding boxes within the 3766 generated Schematismus documents have been calculated and used as a guideline to set the anchor-generation parameters. A general characteristic of the anchor boxes chosen is their elongated and narrow shape, which is understandable, given that text lines have a similar shape. Additionally, some



Table 3 Lists of the class types used to generate synthetic Schematismus documents

Table with 2 columns: Class name and Class description. Rows include Paragraph, BigParagraph, H1, H2, H3, H4, NameEntry, and Curly.

1360

Name index section listing names and their corresponding page numbers, such as Becker Alexander, Beckh Rudolf, Beckmann Friedrich, etc.

Fig. 8 Example of an original page from the name index section

Synthetic documents are primarily composed of paragraphs. Paragraphs always begin with a bolded word and are indented after the first line.

A big paragraph is a paragraph with a bold starting word, but an inverted indentation compared to a normal paragraph. In addition, this type of structure element may never appear within a column.

H1 headlines represent headings that appear above the start of a column. This class may also never appear within a column.

H2 headlines have a slightly smaller font size than H1 headings, but will always be contained within a column.

An H3 header is a title that has the same font size as a paragraph and may only appear within a column. The class is also used to define keywords on the right-hand side of a curly bracket.

An H4 header is formatted in the same way as an H3, but italicized and always enclosed in parentheses.

Name entries consist of a bold text that is left aligned followed by one or more numbers that are right aligned.

Using the curly class, multiple paragraphs and a single H3 class can be contained within. In this manner, it should be easier to assign individual paragraphs to a H3 class in the future.

Example of generated Schematismus-style document with annotations. Includes sections like 'Jubiläums-Erinnerungsmedaille für die bewaffnete Macht.', 'Altenburg. I. Seitzthal.', and various name entries with class labels and bounding boxes.

Fig. 9 Example of generated Schematismus-style document alongside with corresponding annotations (bounding boxes and class label) used in the training set

objects, such as single lines or headings, were quite small, so the anchor boxes generated were smaller than usual.

3. Further, the maximum number of object detections per image needed to be adjusted. Since most object detection models detect only a few objects at a time within a single image, the default value is 100 objects. It should be noted, however, that since the purpose of this analysis is to detect quite fine-grained layout elements within these *Schematismus* documents, the number of layout elements within one document may easily exceed 250 (some pages, for instance from the name register of the *Schematismus*, feature up to 400 objects on one page). In order to completely disregard this upper limit of detections per image, the parameter was set to 1000.
4. Another significant adjustment has been made to the resolution of the images. It should be noted that while the standard image input resolution of the faster R-CNN model is  $1333 \times 800$  pixels, this resolution results in the model sometimes being unable to detect smaller objects such as single lines within the *Schematismus* documents. The reason for this is that both convolution layers and pooling layers in the model further downscale the input image, resulting in loss of important information. Based on the experiments conducted, a resolution of  $1988 \times 1405$  pixels has been determined as the input resolution.

As mentioned previously, 3766 synthetic *Schematismus*-style documents have been generated for the purpose of creating a fully annotated data set. Among these, 3126 served as the training set and 640 served as the validation set. A stratified split of the full data set was used to select the training and validation sets. In our scenario, the stratification process involved counting the occurrences of each class within each generated document page. This information was then used to split the data into two sets, maintaining a similar class distribution in both the validation and training sets. Even though there were more than 3000 annotated training documents available, dataset augmentation strategies have been applied. Adding random augmentations to existing data allowed the training set to be artificially expanded in terms of document variety without increasing the number of documents and thus increasing training time. Therefore, adding augmentations to the training of a faster R-CNN layout detection model of documents can improve its accuracy and robustness.

### 3.3.2 Training data augmentation and further steps

As a result of applying random transformations such as rotation, scaling, cropping, optical distortion (to simulate page warping), blur, noise adding and page flipping, the model could learn to recognise and locate different layout elements invariant of their angle or size. Considering that all of these parameters were selected randomly, it is extremely unlikely

that two identical documents will be input into the model during training. Additionally, augmentation may help to reduce overfitting, which occurs when a model becomes too specialised in recognising only training examples and performs poorly on new, unknown data. By augmenting the training data, the model is exposed to a wider range of layout variations and becomes more adaptable to new and unseen documents.

As for the actual training process, adjustments have been made to the number of epochs, batch size, and learning rate. During each training iteration, batch size determines the number of samples, and thus images, to be processed by the machine learning model before the weights are updated. It is one of the most influential hyperparameters when training deep learning models, and it can be viewed as a trade-off between accuracy and speed. A larger batch size allows more samples to be processed at once, resulting in faster training times and better hardware utilisation. However, larger batch sizes require more memory and may hinder generalisation [45].

By contrast, a smaller batch size results in fewer samples being processed at once. Despite slower training times, this can also prevent, to some extent, overfitting and produce a more generalisable model. Typically, smaller batch sizes are used when the data set is small or when the model requires frequent updating of a large number of parameters. Choosing the correct batch size cannot be achieved in a one-size-fits-all manner since it is heavily dependent on the data set being used. For our study, a batch-size of two has been selected based on experimentation since it fully utilises GPU memory and, along with a scaling factor of 85 percent, produces a relatively fast training process.

In order to maintain a constant variance in gradient expectations, it is recommended to multiply the learning rate by  $\sqrt{k}$  when multiplying the batch size by  $k$  [46]. As a result of extensive learning rate optimisation, a base learning rate of 0.005 has been chosen for batch size one. Thus, the final learning rate is  $0.005 \cdot \sqrt{2} \approx 0.007$ . This learning rate, along with a weight decay of 0.0005 and a momentum of 0.9, was used to initialise a stochastic gradient descent optimiser. The total number of epochs was set to 100. The model is saved after every epoch if the validation loss is less than the previous saved validation loss. In this manner, one can be assured that the final model, which has been trained for 100 epochs, is the one which worked best on the validation set.

The model, which has been trained purely on synthetic data, gave fairly solid results when used on original documents, as described in detail in the [Evaluation](#) section. While these results are already promising, the existing model can also be further fine-tuned using real, original documents. As annotations could not be generated this time, they had to be hand drawn, which is a very time-consuming process.

PyTorch's TorchServe [47] framework is, however, accelerated this process significantly. Employing this method allowed us to "serve" an existing trained model over the local network, where one could send images to and receive bounding-box and label prediction information. In theory, this is not much different from a simple Python script that runs a model directly to predict the layout elements of documents that are fed into it, but there are some very specific applications it can be used for. With the help of an annotation software called "BoundingBoxEditor" [48] the pre-trained model can be "served" and therefore accelerate the manual annotation process of the original documents significantly. This is due to the fact that the pre-trained model already yielded good results. Therefore, only a few adjustments and error corrections were necessary, such as correcting incorrect classifications or bounding boxes. A total of 39 original *Schematismus* documents have been manually annotated and saved in this way. Once the original *Schematismus* documents had been annotated, the annotations were utilised to fine-tune the existing model. To that purpose, the newly created data set had been divided into training and validation sets. Using the pre-trained model's weights as initial weights, these sets were trained for 100 epochs using the same parameters described earlier. Figure 12 show the training loss and validation loss for this training process.

### 3.3.3 Layout detection post-processing

When analysing the predictions in detail, it becomes apparent that some bounding boxes are overlapped, leading to sometimes incorrect predictions. An illustration of this phenomenon can be found in Figure 11, in which two different bounding boxes overlap on the second line.

To address this issue, a bounding box and label classification post-processing step has been developed. The process works by iterating over every predicted bounding box and then calculating the intersection over union (IoU) with every other bounding box. Essentially, the IoU represents the ratio between the overlapping area and the union area, and the closer the IoU is to 1.0, the more similar the bounding boxes are. As soon as the IoU has been calculated for every bounding box, merge candidates are identified by selecting boxes with an IoU score higher than 0.3. As a result, a maximum bounding box is calculated, which encompasses all merge candidates (including the box currently being viewed), and is used to replace the original bounding box. Considering that the merge candidates may be of different classes, the merged bounding box will be labelled with the class tag of the highest confidence score. It should be noted that bounding boxes associated with the "Curly" class will not be selected as merge candidates, as the nature of this class is to have enclosed boxes inside them. Figure 11 illustrates the over-

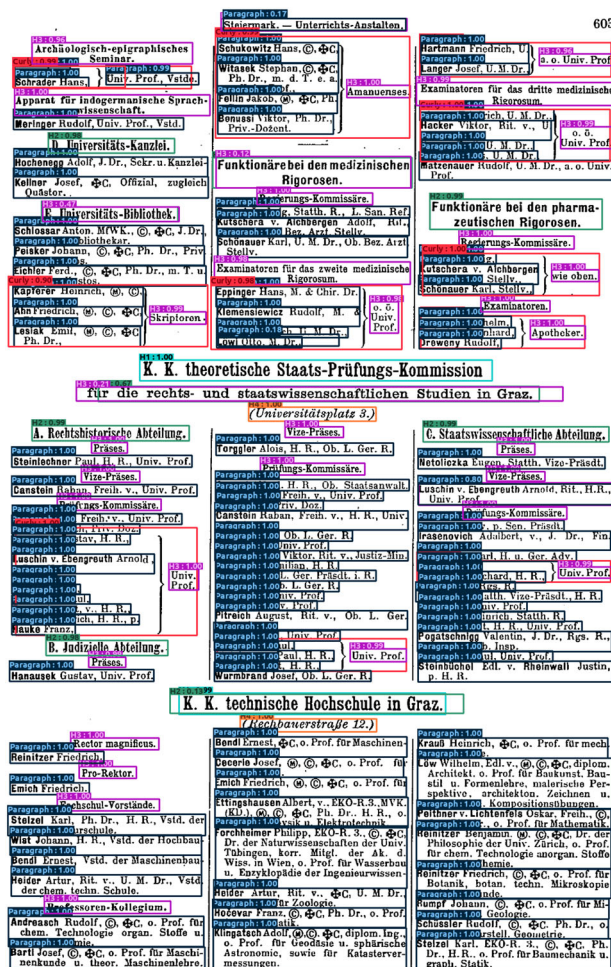


Fig. 10 This figure illustrates a randomly selected original *Schematismus* document with layout elements overlaid. Element predictions with confidence levels less than 0.1 have been omitted. The faster R-CNN model used to get these prediction results has been fine-tuned on 39 original *Schematismus* documents



Fig. 11 The figure illustrates that some bounding boxes overlap and have ambiguous label classifications, there are two overlapping bounding boxes in the middle (magenta and teal)

lay of the predicted boxes following the application of the post-processing step.

### 3.4 OCR

To extract the text within the individual elements of the predicted layout, Tesseract 5.0 was used [49]. As mentioned in Subsection 3.2, the font used in *Schematismus* documents is no longer commonly used. Despite the fact that Tesseract has been pre-trained on a number of fonts, it makes sense to utilise the custom fonts developed for the generation of

*Schematismus* documents to fine-tune the Tesseract optical character recognition model. In addition, due to the symbols that are used in the documents (see Fig. 4) which are unique to *Schematismus* documents, training on this font is necessary in order to recognise these symbols. To fine-tune Tesseract on such a font, images containing a single block of text rendered in the font must be generated. For this, Tesseract's built-in function "text2image" has been utilised. This function requires a large textfile of training text as one of its parameters. Despite the fact that there are existing text files for fine-tuning in German, a custom text file has been compiled using the same text generation method as described in Sect. 3.2.

In addition, a custom character mapping file called "unicharset" must be provided in order to map the *Schematismus* unique symbols to special Unicode characters. The built-in function has been used to create 50,000 images in total. Additionally, to generating individual images containing a single text block and saving them as a ".tif" file, two additional files are generated. One of them is the underlying ground truth, which is saved as a text file. As for the second file, it contains information about every character rendered within the image. It provides details about the character as well as coordinates describing its bounding box. The three files are then combined into a single ".lstmf" file, which is essential for the training process once fine-tuning begins. A note should be made regarding the fact that the German OCR model has been used as a starting point for this training process. Section 4 presents an evaluation of the fine-tuned Tesseract model, as well as additional experiments and pre-processing steps required to obtain the best results.

## 4 Evaluation

We start the evaluation of our model performance with an explanation of our parameter choices for the layout detection model, discussed in subsection [Layout detection](#). Then, both the Tesseract optical character recognition algorithm and the layout detection model are evaluated in detail. As these two elements perform quite distinctly, evaluation will first take place separately, followed by an evaluation of the combined results in subsection 4.3. Furthermore, any additional pre- or post-processing steps that could be performed to further improve the results will be described.

### 4.1 Evaluation of the layout detection model

A set of eighteen original *Schematismus* document pages from 1910 was analysed for evaluation of the layout detection model. All the original documents are completely new to the model. It is important to note that these documents were not part of the training or validation set used to fine-tune

the model. An example of how predicted bounding boxes and corresponding labels appear is shown in Figure 10. This figure illustrates the model being applied to one of eighteen selected document pages. Predicted elements feature a confidence level between 0 and 1, representing the model's certainty about the accuracy of its prediction. Boxes, which feature a confidence level below 0.1, have been omitted.

Overlapping bounding boxes were sorted out, using the post-processing step described in subsection 3.3.3. As a result of this step, the layout detection model can now be evaluated. For the purpose of obtaining a ground-truth of bounding-boxes with corresponding labels, the eighteen selected original documents are hand-annotated. During the drawing of the bounding boxes, extensive attention has been paid to detail, in order to reproduce a very accurate ground truth.

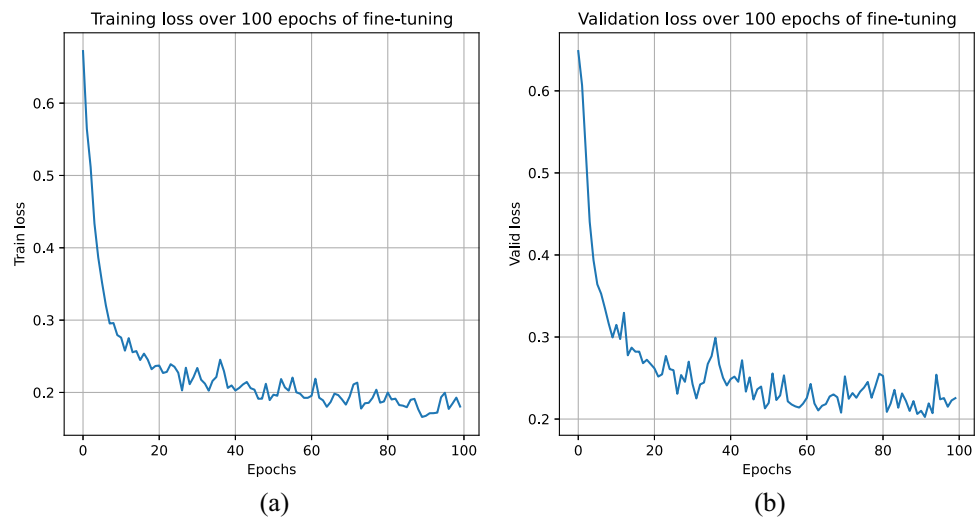
To measure the predicted bounding box accuracy, the intersection over union (IoU) method is again employed. For each document image in the testing-set, the bounding boxes in the ground-truth set were iterated over and the best matching prediction based on the IoU score was selected. Following the establishment of a list of all best matching predictions, the average over all IoU scores is calculated, which represents the bounding box prediction accuracy for the given page.

It should be noted that in order to distinguish between binary values in a multiclass classification problem, the metrics have to be calculated for each class individually. This is while only considering the class currently in focus to be positive (1) and all the rest as negative (0). Figure 13 illustrates this with a confusion matrix. In order to measure the performance of the classification of each layout element, four different metrics are used: accuracy, precision, recall, F<sub>1</sub>-score (Fig. 14).

All the metrics mentioned above have been calculated for each of the eighteen documents in the test set in accordance with the ground truth. Table 4 gives a detailed overview of each metric for both fine-tuned and non-fine-tuned models. Compared to the other documents, the documents with indices 3 and 5 performed the worst, especially in terms of accuracy. Explanations and illustrations of why these performed so poorly are provided in Sect. 5.1. Apart from that, the results for both bounding box accuracy and classification performance appear promising. Table 5 lists the final performance metrics of the fine-tuned model, averaged across all tested pages.

Additionally, Table 6 provides a statistic about the confidence with which the faster R-CNN model has predicted bounding-box and corresponding class labels for each class. There are three approaches that can be used in order to gain insight into this behavior: The first column in the table represents the average confidence level whenever anything has been detected by the model. Column two indicates the confidence level of the model's prediction, when the associated

**Fig. 12** This figure illustrates the training and validation losses associated with a faster R-CNN model trained on a dataset containing 39 original documents. Initial weights were derived from the weights of the existing pre-trained model



**Table 4** Summary of classification accuracy, precision, recall, F<sub>1</sub>-score, and bounding-box accuracy (based on the average IoU) for both non-fine-tuned and fine-tuned models

Document- index	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
<i>Fine-tuned</i>																		
Accuracy	0.98	1	1	0.74	1	0.8	1	1	1	0.97	0.95	1	0.99	1	1	1	0.99	1
Precision	0.99	1	1	0.74	1	0.56	1	1	1	0.98	0.75	1	0.97	1	1	1	0.73	1
Recall	0.94	1	1	0.75	1	0.69	1	1	1	0.92	0.74	1	1	1	1	1	0.62	1
F <sub>1</sub>	0.95	1	1	0.72	1	0.57	1	1	1	0.95	0.71	1	0.98	1	1	1	0.66	1
bbox- Accuracy	0.91	0.92	0.92	0.90	0.89	0.89	0.79	0.80	0.79	0.96	0.95	0.96	0.97	0.95	0.97	0.97	0.85	0.97
<i>Non fine-tuned</i>																		
Accuracy	0.81	0.89	0.67	0.81	0.93	0.50	0.97	0.99	0.97	0.83	0.78	0.66	0.59	0.69	0.52	0.54	0.88	0.57
Precision	0.51	0.75	0.38	0.49	0.65	0.50	0.50	0.50	0.50	0.47	0.47	0.64	0.71	0.65	0.55	0.66	0.41	0.63
Recall	0.54	0.71	0.32	0.51	0.53	0.38	0.49	0.49	0.49	0.44	0.43	0.45	0.56	0.55	0.50	0.44	0.45	0.51
F <sub>1</sub>	0.52	0.71	0.35	0.49	0.56	0.42	0.49	0.50	0.49	0.44	0.45	0.52	0.62	0.59	0.48	0.53	0.43	0.55
bbox- Accuracy	0.77	0.78	0.72	0.70	0.73	0.71	0.75	0.76	0.75	0.79	0.78	0.79	0.73	0.71	0.76	0.71	0.79	0.68

**Table 5** Performance measures were calculated based on the average of all 18 test documents

Performance measure	Value
Average classification accuracy	0.968
Average classification precision	0.929
Average classification recall	0.925
Average classification F <sub>1</sub>	0.919
Average bounding-box accuracy	0.909

class was actually correct, and column three indicates the confidence level when the associated class was incorrect. Whenever no incorrect predictions have been made for a particular class, the value has been omitted. Based on the results in column three, it can be seen that the model tends to be overconfident in its predictions.

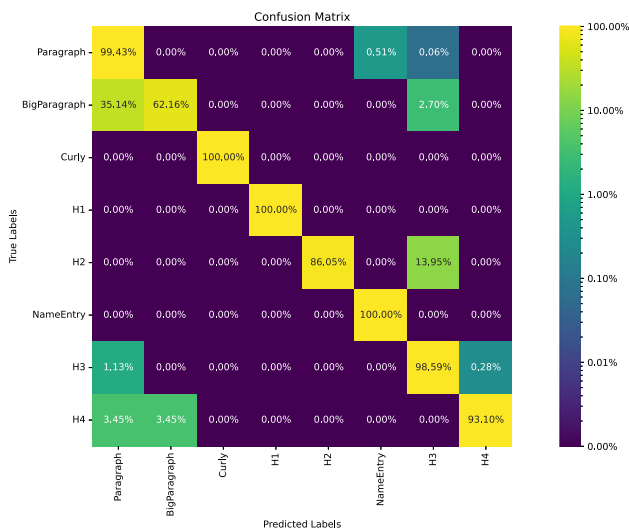
### 4.2 Evaluation of OCR performance

To see how Tesseract OCR performs on original *Schematismus* documents, a ground truth must be established. As this ground truth must be compiled manually by converting documents into plain text, this takes a considerable amount of time. For the evaluation process, a total of 16 original *Schematismus* pages have been manually transcribed this way. Then, in a first step, all pages were fed into a distribution of Tesseract OCR, which had not been fine-tuned to the custom font, furthermore, the pages were not preprocessed via layout detection. Tesseract was configured to use built-in page segmentation to partition the outputs of the entire pages into an easily readable and correct format. The resulting outputs did not match expected sequence. Consequently, in order to make a fair comparison between Tesseract’s output and the corresponding ground-truth, the predicted outputs have been manually split and reordered to match the original

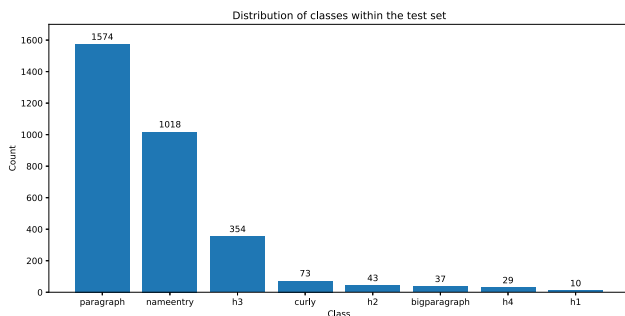
**Table 6** In the Table the confidence levels of the faster R-CNN model in predicting bounding-boxes and corresponding class labels for each class are displayed

Class name	Average confidence for any prediction	Average confidence for correct prediction	Average confidence for incorrect prediction
Paragraph	0.983	0.991	0.899
BigParagraph	0.923	0.924	<b>0.896</b>
Curly	0.970	0.970	–
H1	<b>0.997</b>	<b>0.997</b>	–
H2	0.911	0.911	–
H3	0.961	0.988	0.797
H4	0.958	0.975	0.747
NameEntry	0.749	0.967	0.422

Presented are the average confidence level, the confidence level when the associated class was correct, and the confidence level when the associated class was incorrect. According to the results, the model tends to overestimate its predictions, as evidenced by the higher confidence levels associated with incorrect predictions



**Fig. 13** This figure illustrates the confusion matrix for all predictions made on the test set. Note that the values have been normalized according to the ground truth, so that each row sums to 100 percent



**Fig. 14** Using ground truth, this figure illustrates the distribution of the various layout elements found in the test set

layout of the specific page, without altering any extracted characters or words. The constrained number of tested pages is also attributed to the time-intensive nature of this additional process. The exact same process was then repeated using a distribution of Tesseract OCR, which had been fine-tuned for the custom font. Following manual alignment, CER and WER were calculated for every block of text. The average over all pages is shown in rows one and two of Table 7. According to the results, the average CER has improved by 7.01 percentage points and the average WER by 10.94 percentage points when using the fine-tuned OCR model.

### 4.3 Evaluation of the OCR in combination with the layout detection model

In the next step, the layout detection model was utilised to segment every structure element within the pages into individual elements. To obtain the image snippets, the original images were cropped based on the predicted bounding boxes for each page. As each element’s coordinates are known in the original document, it was possible to sort them in the appropriate order, so no manual reordering was necessary. The Tesseract OCR model was then applied to each image snippet individually. Since it became evident in the prior step that the fine-tuned version performs significantly better, only this version was used. For each individual image snippet, CER and WER were calculated with the corresponding ground truth. Row three in Table 7, the averages of all predictions across all pages are presented.

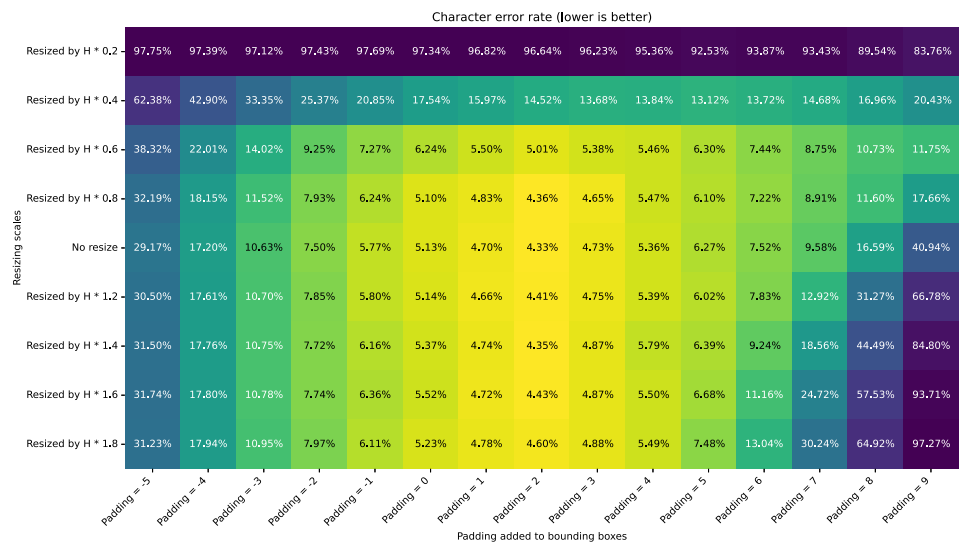
According to these results, both average CER and WER have improved by another 7.87 and 6.38 percentage points respectively compared to the averages computed based on feeding the full pages into the fine-tuned Tesseract OCR algorithm. Even though we consider these improvements satisfying, we believe that OCR accuracy can be further improved. Although the layout detection model performs

**Table 7** The table presents the average CER and WER resulting from various scenarios

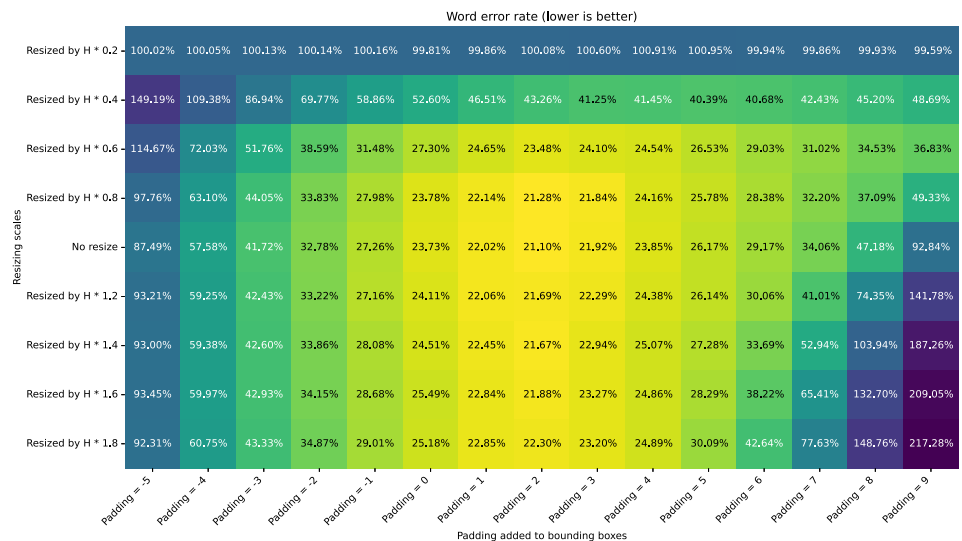
	Avg. CER (%)	Avg. WER (%)
Full page-without font fine-tuning	20.01	41.05
Full page-with font fine-tuning	13.00	30.11
Utilizing layout detection	5.13	23.73
Final score	4.33	21.10

The first two rows show the impact of fine-tuning Tesseract to a custom font on CER and WER on a full page. The third row presents the results when layout detection is combined with a fine-tuned Tesseract OCR model. Finally, in the last row, the results after a padding of two pixels and no resizing are applied to each snipped image, are shown

**Fig. 15** The figure shows the average character-error-rate and word-error-rate calculated with different levels of upscaling/downscaling and padding applied to individual image snippets



(a) average character-error-rate



(b) average word-error-rate

well in finding very accurate bounding boxes, sometimes characters are cropped off at the borders of the images. As a result, padding has been added around the predicted bounding box, to address this issue. As described in the official Tesseract guide on improving OCR accuracy [50], Tesseract generally works better with higher resolution images. Therefore, the individual image snippets have been resized using various scales based on the height of each image in order to find the sweet spot. The aspect ratio of the original image was maintained while upscaling (or downscaling) through linear interpolation. Figure 15 illustrate this. Interestingly, the CER and WER reach their minimum values when the original image snippet is left unscanned and padded with two pixels. This observation suggests that the resolution of the scanned pages may be sufficiently high. Row four in Table 7 shows the final average CER and WER values.

According to these results, it is possible to answer the research question, which is how much OCR accuracy can be improved by using a layout detection model as a preprocessing step to segment *Schematismus*-state documents, and feeding Tesseract individual images containing one layout structure rather than a full page. In comparison to the average CER and WER obtained on a full page using a fine-tuned Tesseract OCR model, an 8.67 percentage point improvement in the CER and 9.01 percentage point improvement in the WER were observed. In comparison with an out-of-the-box Tesseract OCR model, even higher CER and WER improvements were achieved. Here the average CER improved by a total of 15.68 percentage points and the average WER improved by a total of 19.95 percentage points.

### 5 Discussion

It is evident from the results presented in section 4 that the use of a custom-developed layout detection model to segment *Schematismus*-style documents together with a Tesseract model fine-tuned to a custom font designed to be as close to the original as possible significantly improved the quality of the extracted text. However, it should be noted that, due to the relatively small sample size of 16 pages, these results might not represent the full picture. The gain that combined layout detection and text extraction provides may be larger than metrics alone can express. That is due to two different reasons.

As each layout element is segmented by bounding boxes, it is known where the blocks are located within a document coordinate-wise. It allows the reordering of the extracted texts so that they correspond to the reading flow of the document. This is essential when extracting text from a document with columns.

The second reason is that due to the classification of the individual bounding boxes it is possible to immediately tell

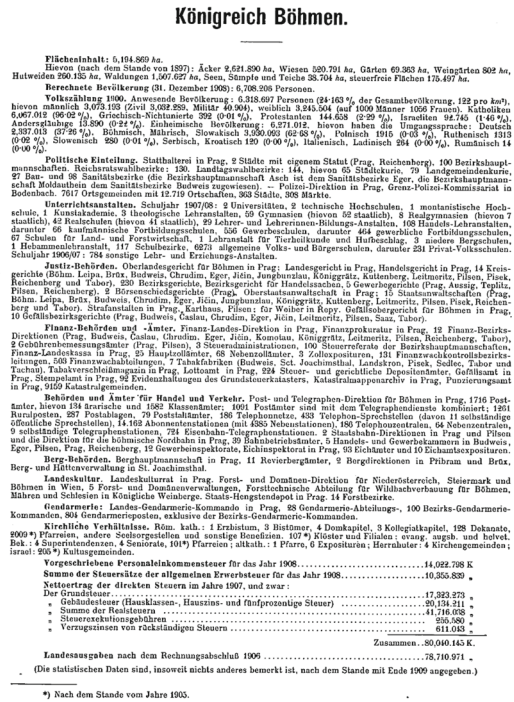


Fig. 16 This figure illustrates a document with index 5 of the evaluation set, which had the poorest layout-detection performance

the class of a structure element. Thus, it is possible, for example, to extract only headings and paragraphs from a document. Moreover, this makes it easier to match paragraphs with individual headings, or to get all enclosed structure elements within curly brackets.

As for the Tesseract OCR model, the custom designed font should be improved, to make text extraction even better. Due to the fact that the current version of the font does not include certain characters such as “č” or “ň” in its unichar set, detection of these kinds of letters is not possible, resulting in errors.

The pipeline we outlined in this article is highly adaptable and could be optimised for similar tasks with relatively little effort. Particularly the process we suggest to produce synthetic training data at scale contributes largely to our capacity to adapt the layout detection model quickly. Not only does it allow to produce a significant amount of training data in relatively little time, but it also does of synthetic training data also appears to significantly increase the precision of the bounding box prediction.



## 5.1 Error analysis

According to Table 4, documents with the indices 3 and 5 performed quite poorly compared to others in terms of accuracy. Both documents are visually very similar. As an example, the document with index 5 is shown in Fig. 16. Clearly, this document differs from the typical three-column *Schematismus*-style document. Aside from the general layout, a key difference is the indentation of each paragraph. Although these aspects were considered during the generation of synthetic documents, resulting in a separate class “BigParagraphs”, the generated structures do not appear to be as similar to the originals as intended, based on layout detection results. Due to the relatively small number of examples of this type of document in the training set used for fine-tuning, we could not observe any improvement. Specifically, only two pages containing “BigParagraphs” were included in the fine-tuning training set, which appears to be too few for the model to effectively learn this class. Therefore, to improve performance on these types of *Schematismus* documents, more pages similar to those in Fig. 16 must be manually annotated and added to the fine-tuning training set.

## 5.2 Research questions

With regard to our first research question, we could show how OCR accuracy can be significantly improved by splitting individual document pages into their layout elements as a preprocessing step.

As for the second research question, it has been shown that fine-tuning Tesseract with a custom font results in performance improvement.

In comparison with the performance of an out-of-the-box Tesseract Model for OCR on an entire page of the *Schematismus*, the results indicated that segmenting and splitting individual document pages into their layout elements with a deep learning convolutional neural network resulted in significantly better OCR accuracy.

## 5.3 Outlook

The procedure we developed therefore represents a crucial step toward a significantly improved analysis of printed historical documents produced in the larger context of the long 19<sup>th</sup> century, particularly as we show how each of the two steps can be further adapted to the specific needs, requirements and challenges met by fellow researchers.

However, we expect that both, layout detection and optical character recognition, can be further optimised for even better performance on historical documents. Layout detection may benefit from increasing the training dataset, not only in terms of the number of pages but also by including a wider variety of documents. In other words, by generating documents that

are visually similar to much older, in our case *Schematismus*-style documents produced in the first half of the 19<sup>th</sup> century, when a different layout was used, and including these in the training set, a more generalized and robust model may be achieved. Further, domain knowledge can be put to use to enhance text extraction. Considering that most of the printed text in *Schematismus*-style documents consists of abbreviations that are listed and described on specific pages within these documents, this information can be utilised to build a custom spell-checking algorithm to correct errors in the text extraction process. Additionally, it would be of interest to explore whether the methods used in this paper can be applied to other types of historical documents.

Our approach offers a viable solution to a number of common problems in dealing with retro-digitised historical texts in historical and humanities research contexts, yet also in industrial application. In this work, it was first shown that the breakdown of the OCR problem, and its solution in several sub-steps, is very promising. However, careful work and precise adjustment of the training data are necessary preconditions to obtain excellent performance. Future work based on our approach will tackle more diverse layouts and broader scope of documents types.

**Funding** Open access funding provided by Graz University of Technology.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

1. Liu, X., Gao, F., Zhang, Q., Zhao, H.: (2019), pp. 32–39
2. Noflatscher, H.: (Böhlau, 2004), pp. 59–75
3. Bauer, V.: *Repertorium territorialer Amtskalender und Amtshandbücher im Alten Reich: Adreß-, Hof-, Staatskalender und Staatshandbücher des 18. Jahrhunderts.*, vol. vol. 2 (Klostermann, 1999)
4. ALEX. Staatshandbuch 1910: Schematismus staat (1910). <https://alex.onb.ac.at/cgi-content/alex?aid=schb&datum=1910&page=597&size=45>. Accessed on: 2022-10-21
5. Raphael, L.: *Die Erben von Bloch und Febvre. Annales-Geschichtsschreibung und nouvelle histoire in Frankreich 1945-1980* (Klett-Cotta, Stuttgart, 1994)
6. Jannidis, F., Kohle, H., Rehbein, M.: Digital Humanities Eine Einführung. J. B. Metzler, Stuttgart (2017). <https://doi.org/10.1007/978-3-476-05446-3>
7. Boros, E., Nguyen, N.K., Lejeune, G., Doucet, A.: Assessing the impact of ocr noise on multilingual event detection over digi-

- tised documents. *Int. J. Digit. Libr.* (2022). <https://doi.org/10.1007/s00799-022-00325-2>
8. Wajer, M.B.W.: Internet Archive OCR Stack in 2021. Switching to Open Source Software (2021). <https://ia601807.us.archive.org/35/items/merlijn-wajer-presentation/merlijn-wajer-presentation-ocr.pdf>
  9. Austrian National Library. Austrian Books Online (2023). <https://www.onb.ac.at/en/digital-offers/austrian-books-online>
  10. Kettunen, K., Koistinen, M., Kervinen, J.: Ground truth ocr sample data of finnish historical newspapers and journals in data improvement validation of a re-ocring process. *LIBER Quarterly* **30**(1), 1–20 (2020). <https://doi.org/10.18352/lq.10322>
  11. Staatsbibliothek, B.: Technologien & Softwareentwicklung (2023). <https://www.digitale-sammlungen.de/de/technologien-und-softwareentwicklung>
  12. G. Markus. Issue 13: OCR. EuropeanaTech Insight is a multimedia publication about R&D developments by the EuropeanaTech Community (2019). <https://pro.europeana.eu/page/issue-13-ocr>
  13. Ehmer, J., Mitterauer, M., Thaller, M.: Wiener Datenbank zur Europäischen Familiengeschichte (2023)
  14. Becker, P., Osterkamp, J.: The Emperor's Desk (2018–2021)
  15. Romberg, M.: The Viennese Court (2020–2023)
  16. Popovici, V., Velková, A.: Social Mobility of Elites (2022)
  17. Engl, E.: OCR-D kompakt: Ergebnisse und Stand der Forschung in der Förderinitiative. *Bibliothek Forschung und Praxis* **44**(2), 218–230 (2020). <https://doi.org/10.1515/bfp-2020-0024>
  18. Martínek, J., Lenc, L., Král, P.: Building an efficient OCR system for historical documents with little training data. *Neural Comput. Appl.* **32**(23), 17209–17227 (2020). <https://doi.org/10.1007/s00521-020-04910-x>
  19. Reul, C., Christ, D., Hartelt, A., Balbach, N., Wehner, M., Springmann, U., Wick, C., Grundig, C., Büttner, A., Puppe, F.: Ocr4all—an open-source tool providing a (semi-)automatic ocr workflow for historical printings (2019). <https://doi.org/10.3390/app9224853>. <http://arxiv.org/abs/1909.04032>
  20. Cordell, R.: Machine Learning + Libraries. A Report on the State of the Field (2020). <https://labs.loc.gov/static/labs/work/reports/Cordell-LOC-ML-report.pdf?locl=blogsig>
  21. Gasparini, A., Kautonen, H.: Understanding artificial intelligence in research libraries-extensive literature review. *LIBER Quart. J. Assoc. Eur. Res. Libr.* (2022). <https://doi.org/10.53377/lq.10934>
  22. Teibenbacher, P., Kramer, D., Göderle, W.: An Inventory of Austrian Census Materials, 1857-1910. Final Report. *Mosaic Working Paper* **190**, 25 (2012)
  23. Zechner, A., Knapp, E., Adelsberger, M.: Prices and Wages in Salzburg and Vienna, c. 1450–1850 An Introduction to the Data. *Vierteljahresschrift für Sozial und Wirtschaftsgeschichte* **108**(4), 263–270 (2021). <https://doi.org/10.25162/VSWG-2021-0016>
  24. Bavouzet, J.: in *The Habsburg Civil Service and Beyond: Bureaucracy and Civil Servants from the Vormärz to the Inter-War Years*, ed. by F. Adlgasser, F. Lindström (Verlag der Österreichischen Akademie der Wissenschaften, Vienna, 2019), pp. 167–186. <https://doi.org/10.2307/j.ctvqgxx26b.11>
  25. Wang, J., Liu, C., Jin, L., Tang, G., Zhang, J., Zhang, S., Wang, Q., Wu, Y., Cai, M.: Towards robust visual information extraction in real world: New dataset and novel solution (2021). [www.aaii.org](http://www.aaii.org)
  26. Douzon, T., Duffner, S., Garcia, C., Espinas, J.: Improving information extraction on business documents with specific pre-training tasks. In *International Workshop on Document Analysis Systems*. (Springer Science and Business Media Deutschland GmbH, 2022), pp. 111–125. [https://doi.org/10.1007/978-3-031-06555-2\\_8](https://doi.org/10.1007/978-3-031-06555-2_8)
  27. Carbonell, M., Fornés, A., Villegas, M., Lladós, J.: A neural model for text localization, transcription and named entity recognition in full pages. *Pattern Recogn. Lett.* **136**, 219–227 (2020). <https://doi.org/10.1016/j.patrec.2020.05.001>
  28. Tarride, S., Maarand, M., Boillet, M., McGrath, J., Capel, E., V'ezina, H., Kermorvant, C.: Large-scale genealogical information extraction from handwritten Quebec parish records. *Int. J. Doc. Anal. Recogn. (IJ DAR)* **26**(3), 255–272 (2023). <https://doi.org/10.1007/s10032-023-00427-w>
  29. Monnier, T., Aubry, M.: in *2020 17th International Conference on Frontiers in Handwriting Recognition (ICFHR)*, pp. 91–96. (IEEE, 2020)
  30. Gruber, I., Ircing, P., Neduchal, P., Hruz, M., Hlaváč, M., Zajíc, Z., Švec, J., Bulín, M.: in *International Conference on Speech and Computer*, pp. 166–175, (Springer, 2020)
  31. M. Shen, H. Lei, in *2015 12th International Conference on Fuzzy Systems and Knowledge Discovery (FSKD)*, pp. 1566–1570, (IEEE, 2015)
  32. Lat, A., Jawahar, C.: in *2018 24th International Conference on Pattern Recognition (ICPR)*, pp. 3162–3167 (IEEE, 2018)
  33. Grüning, T., Leifert, G., Strauß, T., Michael, J., Labahn, R.: A two-stage method for text line detection in historical documents. *Int. J. Doc. Anal. Recogn.* **22**(3), 285–302 (2019). <https://doi.org/10.1007/s10032-019-00332-1>. [arXiv:1802.03345](https://arxiv.org/abs/1802.03345)
  34. Büttner, J., Martinetz, J., El-Hajj, H., Valleriani, M.: Cordeep and the sacrobosco dataset: detection of visual elements in historical documents. *J. Imaging* (2022). <https://doi.org/10.3390/jimaging8100285>
  35. Binmakhshen, G.M., Mahmoud, S.A.: Document layout analysis: a comprehensive survey. *ACM Comput. Surv.* (2019). <https://doi.org/10.1145/3355610>
  36. Boillet, M., Kermorvant, C., Paquet, T.: Multiple document datasets pre-training improves text line detection with deep neural networks. *Proceedings-International Conference on Pattern Recognition* pp. 2134–2141 (2020). <https://doi.org/10.1109/ICPR48806.2021.9412447>. [arXiv:2012.14163](https://arxiv.org/abs/2012.14163)
  37. Xu, Y., Li, M., Cui, L., Huang, S., Wei, F., Zhou, M.: Layoutlm: Pre-training of text and layout for document image understanding. In: *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. pp. 1192–1200 (2020). <https://doi.org/10.1145/3394486.3403172>. <https://arxiv.org/abs/1912.13318>
  38. Biswas, S., Riba, P., Lladós, J., Pal, U.: Beyond document object detection: instance-level segmentation of complex layouts. *Int. J. Doc. Anal. Recogn.* **24**(3), 269–281 (2021). <https://doi.org/10.1007/s10032-021-00380-6>
  39. LuaTeX. *Luatex* (2023). <https://www.luatex.org/>
  40. Remy, P.: Name dataset. <https://github.com/philipperemy/name-dataset> (2021)
  41. FontForge. *Fontforge*. <https://fontforge.org/en-US/> (2023). [Accessed 09-Mar-2023]
  42. Wikipedia. Liste der österreichischen Orden und Ehrenzeichen — Wikipedia, the free encyclopedia. <http://de.wikipedia.org/w/index.php?title=Liste%20der%20C3%B6sterreichischen%20Orden%20und%20Ehrenzeichen&oldid=231609032> (2023). [Online; accessed 09-March-2023]
  43. A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, S. Chintala, in *Advances in Neural Information Processing Systems*, vol. 32, ed. by H. Wallach, H. Larochelle, A. Beygelzimer, F. Alché-Buc, E. Fox, R. Garnett (Curran Associates, Inc., 2019). [https://proceedings.neurips.cc/paper\\_files/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf)
  44. Li, Y., Xie, S., Chen, X., Dollar, P., He, K., Girshick, R.: Benchmarking detection transfer learning with vision transformers (2021). <https://doi.org/10.48550/ARXIV.2111.11429>
  45. Keskar, N.S., Mudigere, D., Nocedal, J., Smelyanskiy, M., Tang, P.T.P.: On large-batch training for deep learning: Generalization

- gap and sharp minima. CoRR **abs/1609.04836**. [arXiv:1609.04836](https://arxiv.org/abs/1609.04836). (2016)
46. Krizhevsky, A.: One weird trick for parallelizing convolutional neural networks (2014). <https://doi.org/10.48550/ARXIV.1404.5997>
47. PyTorch. Torchserve (2023). <https://pytorch.org/serve/index.html>
48. Fleischhacker, M.: Boundingboxeditor. <https://github.com/mfl28/BoundingBoxEditor> (2023)
49. Tesseract. Tesseract (2023). <https://github.com/tesseract-ocr/tesseract>
50. Tesseract. Improving the quality of the output (2022). <https://tesseract-ocr.github.io/tessdoc/ImproveQuality.html>. Accessed on: 2022-10-21

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.