# OCR for Greek polytonic (multi accent) historical printed documents: development, optimization and quality control

author_block">
Anna-Maria Sichani
Sussex Humanities Lab, University of Sussex
East Sussex BN1 9RG, United Kingdom

Panagiotis Kaddas
Computational Intelligence Laboratory, Institute of
Informatics and Telecommunications, National Center
for Scientific Research Demokritos
Agia Paraskevi GR-153 10, Greece

Georgios K. Mikros
Department of Italian Language and Literature School
of Philosophy, National and Kapodistrian University
of Athens
Athens GR-157 84, Greece

Basilis Gatos
Computational Intelligence Laboratory, Institute of
Informatics and Telecommunications, National Center
for Scientific Research Demokritos
Agia Paraskevi GR-153 10, Greece

## ABSTRACT

abstract">
This paper presents the development and implementation of a robust OCR tool and a related comprehensive workflow for the recognition of Greek printed polytonic scripts. This project is initiated and developed by an interdisciplinary team with expertise in the areas of document image processing, character segmentation and recognition, machine learning, corpus creation and digital humanities. Our paper aims to describe the design and development of the workflow around this project, including data gathering and structuring, OCR tool development, user interface development, experiments on the training procedure of the tool, evaluation, post-correction and quality control of the results.

## KEYWORDS

Optical Character Recognition, Greek polytonic scripts, historical printed documents, image processing, page segmentation, machine learning, post correction workflow

publication_info">
**ACM Reference Format:**
Anna-Maria Sichani, Panagiotis Kaddas, Georgios K. Mikros, and Basilis Gatos. 2019. OCR for Greek polytonic (multi accent) historical printed documents: development, optimization and quality control. In *3rd International Conference on Digital Access to Textual Cultural Heritage (DATeCH2019), May 8–10, 2019, Brussels, Belgium.* ACM, New York, NY, USA, 5 pages. https://doi.org/10.1145/3322905.3322926

boilerplate">
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
*DATeCH2019, May 8–10, 2019, Brussels, Belgium*
© 2019 Association for Computing Machinery.
ACM ISBN 978-1-4503-7194-0/19/05...$15.00
https://doi.org/10.1145/3322905.3322926

## 1 INTRODUCTION

Although the accurate recognition of Latin machine-printed text is now considered largely a solved problem, recognition of non-Latin minority languages or scripts having a large number of character classes is still subject of active research. Greek polytonic (multi accent) scripts have a large variety of diacritic marks and as a result a large number of character classes (more than 270). Due to that, Greek polytonic documents (both handwritten and printed) cannot be successfully processed by the currently available Optical Character Recognition (OCR) technologies. Furthermore, as the Greek polytonic scripts were actively used from around 200 BC to modern times until 1982, we can easily observe that a huge amount of digitised (handwritten and printed) Greek documents still remains without full text processing and analysis capabilities.

Drawing expertise from the areas of document image processing, page segmentation, pattern recognition, machine learning and digital humanities, our paper aims to describe the development and implementation of a robust OCR tool and workflow for the recognition of printed Greek polytonic scripts. This interdisciplinary work is partially supported by and initiated within the COST action [1] Distant Reading for European Literary History(CA16204), aiming to create the first Modern Greek Literary Textual Collection as part of a larger European Literary Textual Collection (ELTeC) and to apply Distant Reading methods in order to revisit the European Literary History. The interdisciplinary team has been working on this project for almost a year now on a semi-funded basis and is currently formed by 7 members from the areas of Digital Humanities, Computational Linguistics and Corpus creation, Modern Greek Literature and Computer Science.

In what follows, we are briefly discussing the state-of-the-art and previous work in the field of OCR for Greek polytonic scripts (Section 2). Then we are describing our workflow, currently consisted of four (4) stages. Firstly, we will discuss the data acquisition process and data quality specifications, while

---

[1] https://www.distant-reading.net

footer_navigation">9

discussing obstacles and common problems in data gathering and corpus creation (Section 3). Furthermore, we will present the development of a tool for the recognition of old Greek polytonic scripts, by using very simple machine learning algorithms as well as the development of a ground-truthing Graphical User Interface (GUI) application for the manual annotation of old Greek polytonic documents (Section 4). We will then focus on the application and the evaluation of the trained model for the recognition of the data corpus (Section 5). Lastly, we will discuss the development and implementation of a quality control and post-correction framework for the resulted OCR-ed documents (Section 6). We will conclude with some plans for future development of this research (Section 7).

## 2 PREVIOUS WORK

The automated recognition of Greek printed polytonic scripts was for many years an open question/challenge for researchers in the Humanities and Computer Science.

Initiated by the community of (Digital) Classics, OCR for Greek polytonic (mainly Ancient Greek) is still a claim and an open challenge for Humanities and Cultural Heritage researchers working with textual resources using Greek polytonic. Several previous attempts developing an OCR engine [2], focusing on line recognition and character patterns, alongside a GUI interface [3], are offering limited accuracy and poor results for non-Unicode (old) typefaces [4].

Related work within the Perseus project in [3] on progressive multiple alignment and constrained spell checking to different OCR outputs of printed critical editions of Ancient Greek texts is a promising development in Greek polytonic OCR.

In [6], an OCR framework for the recognition of machine-printed Greek polytonic documents is proposed that is based on combining different recognition modules. One module is used for accent recognition while four modules are used for recognizing characters that belong to different horizontal text zones. Moreover, pre-processing stages like text dewarping and text line detection are included in order to assist the recognition modules.

In [5], adaptive zoning features are proposed as an improvement to the already robust zoning feature descriptors. The main idea is to adjust the position of every zone based on local pattern information. Movements are based on the maximization of the local pixel density around each zone.

In [7], a recognition system for ancient Greek polytonic scripts is based on Hidden Markov Models (HMM) [1] is used. This system is character segmentation free and can be applied both on character and word level. The HMM-based method uses language models of bigrams at character level and yields more accurate results when compared to previous techniques which require a more sophisticated character segmentation technique.

---

[2]http://tesseract.projectnaptha.com
[3]https://dcthree.github.io/antigrapheus/
[4]https://unicode.org

## 3 DATA ACQUISITION

In order to train the OCR tool, a corpus of Greek texts using polytonic scripts was created as a test bed. As this project was related to and partially funded by the Cost action Distant Reading, the Modern Greek corpus was mainly compiled by hundred (100) 19th century Greek novels and their length ranges between 7-400 pages. There have been also several other criteria for the corpus sampling, in order to create a text corpus that would be balanced in terms of authors'gender representation, date of publication, length, et.c.

The majority of the texts were initially hold and digitised by various Greek institutions, including the National Library of Greece, using different digitisation standards, thus, resulting to variations in terms of quality of digitisation. In addition, as the texts were derived mainly from original printed editions of the 19th and 20th century using various typefaces and diacritics, there were also major challenges in terms of image quality and variations in the digitisation results.

## 4 OCR TOOL
### 4.1 Proposed OCR Method

The proposed system for the recognition of old Greek printed polytonic scripts consists of four very simple modules and can be seen in Figure 1.

The first one is the initial character segmentation and recognition module, where all input document images are fed into an existing OCR tool [2]. The acquired information is used as an initial stage for our system. Next, these results are post-processed manually (second module) in order to improve the quality of the initial character recognition and create an accurate training set.

The third module is a feature extractor algorithm, were each already segmented character is normalized into a fixed size and then splitted into 3x3 zones. Then, from each sub-image pixel intensity is extracted as image feature and the final descriptor for each character is the zone feature concatenation. This module is used in order to create a train set containing samples from all unique character classes. Also and during test phase, the same features are extracted from each new character and used as input for the classification module.

The last module uses the well-known machine learning algorithm, k-NN [4], in order to classify characters to the nearest class of the train set. The decision relies upon the feature values extracted from the previous module. Experiments with different values of the k parameter (k=1, 3, 5, 7) were conducted and k-NN classification algorithm with k=5 gave was the most accurate results (Section 5).

### 4.2 User Interface

A Graphical User Interface (GUI) was created in order to assist the manual annotation of old Greek polytonic documents. This GUI application is easy-to-use and provides a very fast and accurate way of marking each character of a document
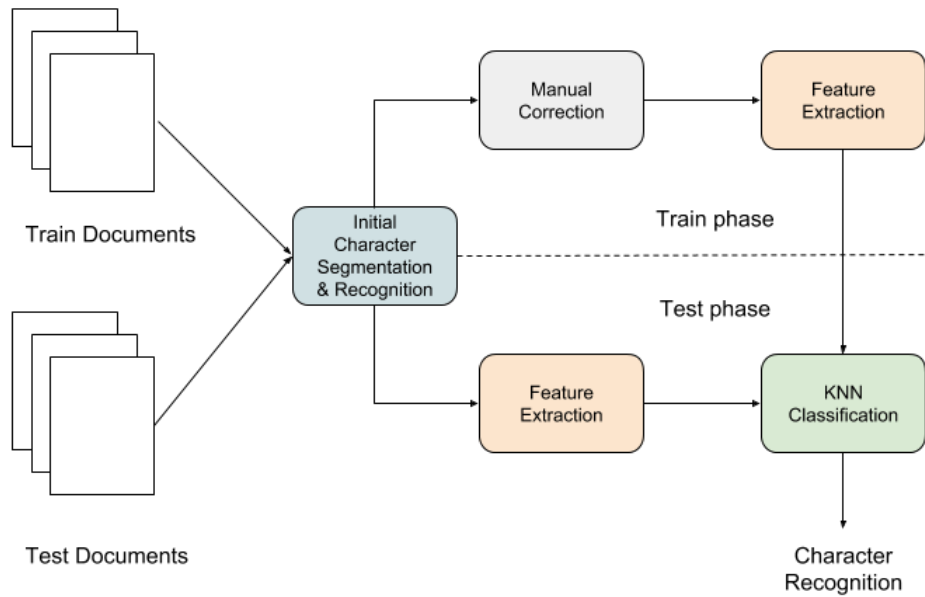
**Figure 1: Overview of the proposed OCR system. Blue: Initial Segmentation & Recognition module used as input to the system. Gray: Manual correction module for input enhancement. Orange: Feature Extraction using pixel intensity for different zones of each character. Green: Classification module using the k-NN algorithm.**
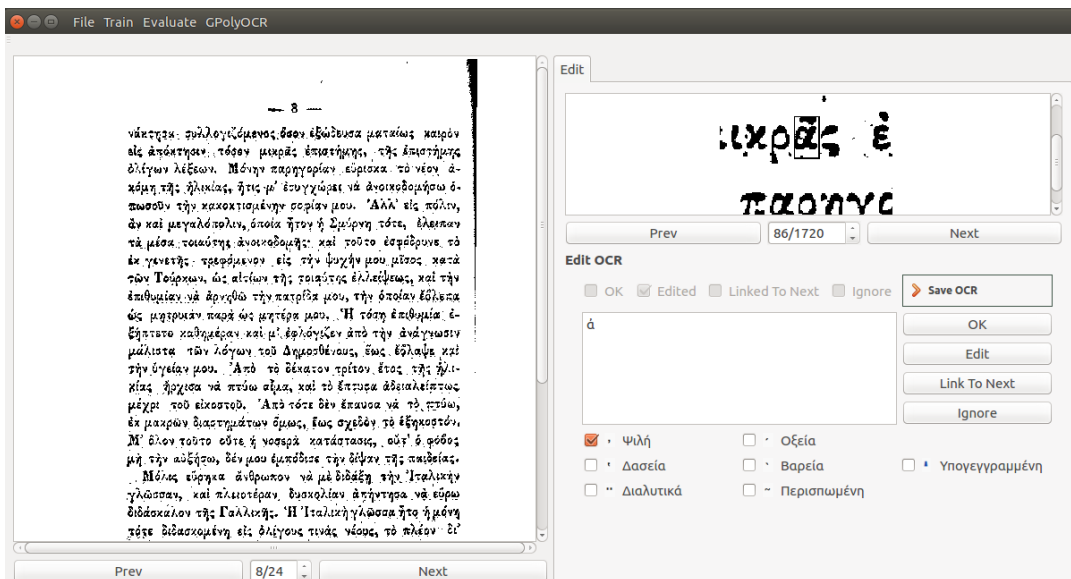


**Figure 2: Example of the developed Graphical User Interface**

page. An existing OCR tool [2] was used in order to acquire an initial character recognition, but recognition results are poor for old Greek polytonic fonts (accuracy is below 80% in most cases) and need further optimization. User feedback can

be embedded by complementing possible polytonic accents to a character just by clicking buttons and without the need of typing special Unicode characters, usually hard to type and impossible to remember. Furthermore, the user is able to discard or modify erroneous character segmentation and recognition results and finally save all changes into a database file. An example of the developed interface is shown in Figure 2.

## 5 EXPERIMENTS

Having developed all the required tools, the training procedure of the model took place by annotating a 1-2 pages per document (depending on the document degradation) using the already described ground-truthing framework. This led into a database system containing information from more than 100 document pages (50k of character entries belonging to 143 unique classes). After completing the training of the model, the OCR tool was applied in batch mode for all the 100 old Greek polytonic novels.

Evaluation of the results using the described framework and random pages as test set yielded accuracy of 94.2% for the best parameter of the k-NN algorithm (k=5) Table 1. This set contained 50k characters belonging to 143 unique classes. The accuracy metric that was used for evaluation was the Levenshtein (or Edit) distance between the predicted transcription of each recognized page and the manually annotated transcription. A maximum accuracy threshold for a page was 98.6% and was usually achieved in document pages of low degradations and noise-free fonts. Respectively, the lowest accuracy was 82.3%.

Moreover, the results are somehow bound to the character segmentation results acquired from the existing OCR tool [2] used as initial stage of the developed system. Despite this, the proposed system results into more accurate results ( 25% improvement) when compared to the initial recognition results.

**Table 1: Character recognition results for the proposed system using the Levenshtein/Edit Distance.**

| Method | Accuracy (%) |
| --- | --- |
| ABBYY FineReader [2] | 68.7 |
| k-NN(k=1) | 89.3 |
| k-NN(k=3) | 92.1 |
| **KNN(k=5)** | **94.2** |
| k-NN(k=7) | 93.8 |

## 6 POST-CORRECTION AND QUALITY CONTROL

In order to further assess the resulted OCRed files, a post correction and quality control strategy was designed and put in place. The team working on the post-correction was formed by five Modern Greek Literary scholars with digital humanities skills and experience.

Firstly, the resulted OCRed files (in txt format) were manually compared with the original digitised files and corrected accordingly (Figure 3). When checking a text, the aim was to ensure that the text is accurately compared to the original source, with alterations and corrections made where considered necessary. At this point, an auxiliary resource was also developed as a look-up vocabulary to assist with and speed up the post-correction phase. The corrected texts were then used as ground truth files in txt format: a comparison with the resulted OCR file could be easily performed, resulting to an evaluation of the entire procedure. The texts are further proceeded and converted in XML TEI format, following the guidelines and the schema of the Distant Reading project working group.

In addition, and in order to document both the quality of the tool and the post-correction resources needed, we are collecting data, mainly concerning three variables: the time needed for the correction, the type and the amount of errors found. With this type of data we are able, alongside qualitative comments, to accurately estimate the resources (both human and financial) needed for post-correction as well as to locate the type of errors and further train the tool for future use.
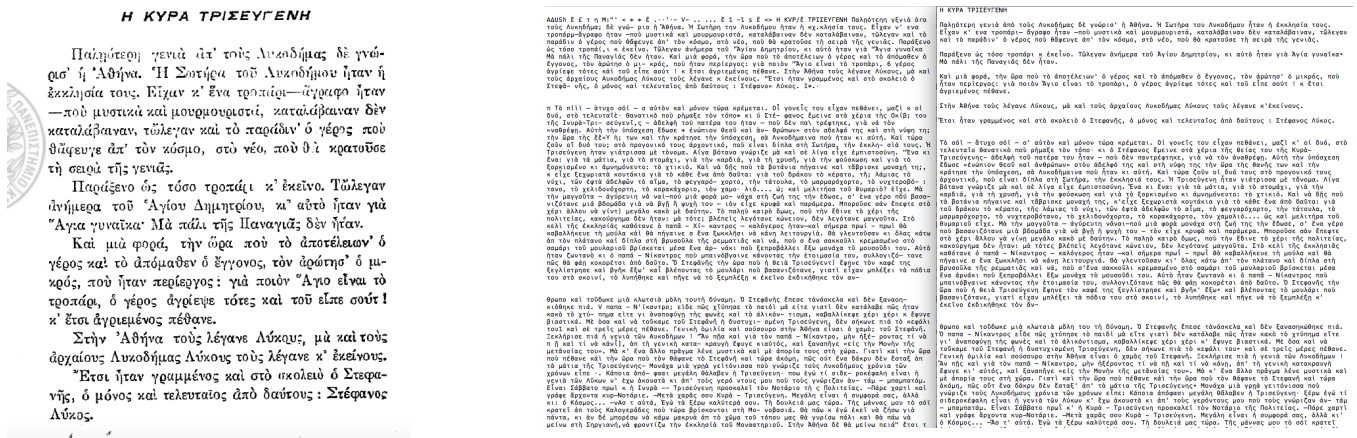
## 7 CONCLUSION & FUTURE PLANS

This research project aims to further contribute to and advance the Optical Character Recognition (OCR) methods for Greek printed polytonic scripts. While working towards the development of the OCR tool, we are also able to point out existing obstacles and gaps in areas such as digitisation standards, ground truth as well as (manual and automated) post-correction. This new OCR framework will facilitate accessing and further processing a huge amount of historical printed polytonic documents, currently digitised and held by numerous cultural institutions but with limited usability in terms of full text access and analysis. By developing a user-friendly and robust OCR framework for Greek polytonic documents, scholars in areas such as Literary Studies, Computational Linguistics, Archival and Information Science will be able to perform previously unavailable advanced computational processes on Greek printed polytonic textual resources.

As for the current phase of the tool, development of a corpus of 19th century printed texts was used as training set; we are willing to enlarge the test set with earlier textual resources from the 18th century as well as latest ones from the 20th century. From a technical point of view, there are a couple of aspects that we are willing to work on in future iterations, such as a more sophisticated classification of characters and perhaps an update on the feature extraction mechanisms.

## 8 ACKNOWLEDGMENTS

**(a) Original Document Image**

**(b) OCR-ed (left) and corrected (right) images**

**Figure 3: Post-correction of the OCR result through comparison with the digitised original source.**

their collaboration in providing access to digitised textual resources.

## REFERENCES

[1] 2008. *Markovian Models for Sequential Data.* Springer London, London, 265–303. https://doi.org/10.1007/978-1-84800-007-0_10

[2] 2019 (accessed January 15, 2019). *ABBYY FineReader.* https://www.abbyy.com/en-us/finereader/

[3] Federico Boschetti, Matteo Romanello, Alison Babeu, David Bamman, and Gregory Crane. 2009. Improving OCR Accuracy for Classical Critical Editions. In *Research and Advanced Technology for Digital Libraries*, Maristella Agosti, José Borbinha, Sarantos Kapidakis, Christos Papatheodorou, and Giannis Tsakonas (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 156–167.

[4] Thomas M. Cover and Peter E. Hart. 1967. Nearest neighbor pattern classification. *IEEE Trans. Information Theory* 13 (1967), 21–27.

[5] B. Gatos, A. L. Kesidis, and A. Papandreou. 2011. Adaptive Zoning Features for Character and Word Recognition. In *2011 International Conference on Document Analysis and Recognition*. 1160–1164. https://doi.org/10.1109/ICDAR.2011.234

[6] B Gatos, Georgios Louloudis, and Nikolaos Stamatopoulos. 2011. Greek Polytonic OCR Based on Efficient Character Class Number Reduction. *Proceedings of the International Conference on Document Analysis and Recognition, ICDAR*, 1155 – 1159. https://doi.org/10.1109/ICDAR.2011.233

[7] V. Katsouros, V. Papavassiliou, F. Simistira, and B. Gatos. 2016. Recognition of Greek Polytonic on Historical Degraded Texts Using HMMs. In *2016 12th IAPR Workshop on Document Analysis Systems (DAS)*. 346–351. https://doi.org/10.1109/DAS.2016.60