

Schriften des Instituts für Dokumentologie und Editorik – Band 3

# **Kodikologie und Paläographie im digitalen Zeitalter 2**

---

## **Codicology and Palaeography in the Digital Age 2**

herausgegeben von | edited by

Franz Fischer, Christiane Fritze, Georg Vogeler

unter Mitarbeit von | in collaboration with

Bernhard Assmann, Malte Rehbein, Patrick Sahle

2010

BoD, Norderstedt

**Bibliografische Information der Deutschen Nationalbibliothek:**

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.d-nb.de/> abrufbar.

© 2011

Online-Fassung

Herstellung und Verlag der Druckfassung: Books on Demand GmbH, Norderstedt 2010

ISBN: 978-3-8423-5032-8

Einbandgestaltung: Johanna Puhl, basierend auf dem Entwurf von Katharina Weber

Satz: Stefanie Mayer und L<sup>A</sup>T<sub>E</sub>X

# Recognizing Degraded Handwritten Characters

Markus Diem, Robert Sablatnig, Melanie Gau, Heinz Miklas

## Abstract

In this paper, Slavonic manuscripts from the 11<sup>th</sup> century written in Glagolitic script are investigated. State-of-the-art *optical character recognition* methods produce poor results for degraded handwritten document images. This is largely due to a lack of suitable results from basic pre-processing steps such as binarization and image segmentation. Therefore, a new, binarization-free approach will be presented that is independent of pre-processing deficiencies. It additionally incorporates local information in order to recognize also fragmented or faded characters. The proposed algorithm consists of two steps: character classification and character localization. Firstly *scale invariant feature transform* features are extracted and classified using *support vector machines*. On this basis interest points are clustered according to their spatial information. Then, characters are localized and eventually recognized by a weighted voting scheme of pre-classified local descriptors. Preliminary results show that the proposed system can handle highly degraded manuscript images with background noise, e.g. stains, tears, and faded characters.

## Zusammenfassung

In diesem Beitrag werden slawische Manuskripte aus dem 11. Jahrhundert analysiert. Herkömmliche *Optical Character Recognition* (OCR) Systeme erzielen schlechte Resultate auf den beschädigten glagolitischen Schriften, da eine korrekte Buchstabensegmentierung nicht möglich ist. Deshalb wird ein segmentierungsfreies OCR-System vorgestellt, welches keiner Vorverarbeitungsschritte bedarf. Da die Klassifikation auf lokaler Information beruht, ist es möglich auch verblasste Buchstaben bzw. Buchstabenfragmente richtig zu erkennen. Das System besteht aus zwei grundlegenden Methoden: Buchstaben-Klassifizierung und Buchstaben-Lokalisierung. Die Klassifizierung basiert auf lokalen, größeninvarianten Merkmalen, die mit Hilfe von *Support Vector Machines* klassifiziert werden. Nach diesem Schritt existieren mehrere gekennzeichnete Merkmalsvektoren pro Buchstabe. Diese werden im zweiten Schritt durch ein *Clustering* Verfahren zusammengefasst, so dass jedem Buchstaben ein finales Klassenetikett zugewiesen werden kann. Die Ergebnisse zeigen, dass auch beschädigte Dokumente mit diesem System automatisch erfasst werden können.

## 1. Introduction

In the digital age Optical Character Recognition (OCR) has been successfully established for automated document analysis of standardized, typeset text. The automatic decipherment of handwriting, however, still poses difficulties for modern and ancient documents likewise. Even more so this applies to damaged or degraded material, the script of which is no longer readable straightforwardly.

In 2007 the interdisciplinary project “Critical Edition of the New Sinaitic Glagolitic Euchology (Sacramentary) Fragments with the Aid of Modern Technologies” of philologists (University of Vienna), computer scientists (image processing group CLV, Vienna University of Technology) and material chemists (Vienna Academy of Fine Arts) was launched to analyse and edit two – later three – valuable Slavonic manuscripts, parchment codices of the Old Church Slavonic canon dating from the 11<sup>th</sup> century: the so-called *Missale Sinaiticum* (Sin. slav. 5/N), a sacramentary fragment consisting of approximately 70 folia written by one main and two minor hands, and a 28-folia fragment, the *Euchologium Sinaiticum pars nova* (Sin. slav. 1/N), part of the famous Sinaitic Euchology discovered in the 19<sup>th</sup> century. They are written in Glagolitic script, which was created in 862/3 by St. Constantine-Cyril for his mission in Great Moravia<sup>1</sup>, and belong to the complex of new findings made in St. Catherine’s Monastery on Mt. Sinai in 1975. Both codices are of a small format of approximately 140x100 mm and are decorated with colour initials and headline highlighting in yellow and green. Unfortunately, especially the Missal shows extensive damages like faded ink, blurring of the ink, staining due to mould or humidity, degradation of the parchment, e.g. chipping, fragmentation and contortion of folia, and the rare phenomenon of chemical conversion of black into white ink. The manuscripts partly contain palimpsest (re-written) folia (for further reference see Miklas 2000).

In the course of the project we have explored several computational approaches to ease codicological and palaeographic investigations, such as layout analysis, semi-automatic character segmentation and feature extraction using graphetic distinctive features – as opposed to this approach –, initial detection and automated puzzling (cf. Kleber and Sablatnig 2009). Furthermore, we have investigated in particular on the description, decipherment and reconstruction of (latent) texts of the manuscripts in question with methods of multi-spectral imaging, image binarisation, document image analysis and image enhancement (Lettnner et al. forthcoming; Miklas et al. 2008). With the combined approaches the readability of the *Missale Sinaiticum* could be enhanced up to 51% (Miklas et al. forthcoming).

Due to the heavily degraded condition of our objects common OCR methods based on robust binarization algorithms did not show satisfying results. Consequently our

---

<sup>1</sup> Designed for liturgical use, the original Glagolica comprised 36 letters functioning also as numerals and fraught with various aspects of theological symbolism (Miklas 2003)

new system for OCR and character supplementation is based on an entirely binarization free method. It performs three major steps, which will be discussed in the following sections: First, the local features of a whole manuscript page are computed and classified by means of Support Vector Machines (SVM). Then, a clustering of interest points according to their spatial coordinates and scale enables the localization of characters. This results in probabilities for character classes that are the basis of a voting scheme for character labels. Preliminary results show that the proposed system can handle images of severely damaged manuscripts with low contrast of text and background and fragmented characters.

## 2. Related Work – Technical Overview

In this section, state-of-the-art OCR systems for degraded documents are presented. It is not intended to give a comprehensive overview (which was already done by Plamondon and Srihari; Vinciarelli), but to describe current developments in the recognition of historic manuscripts. To our knowledge no OCR system has been proposed that can extract features from gray-scale or color images. Current OCR systems have three basic steps in common: First, a document pre-processing is performed. There, the document's skew is estimated, the text layout extracted, and the document image binarized. Subsequently, binary features are extracted and classified by means of a Neural Network (NN) or a SVM. The approaches differ according to the investigated data. Generally, two data sets can be distinguished: For cursive handwritten documents a word based approach is chosen, for non-cursive, usually ancient manuscripts, a character based approach.

**Cursive handwritten documents:** Lavrenko et al. directly recognize words from documents of the George Washington collection. Their technique was later improved by Rath and Manmatha, who added compensation to non-linear variations present in manuscripts. Another word recognition system is proposed by Frinken and Bunke. They compute statistical moments from sliding windows that are applied to normalized word images. Hofmeister et al. compare sample word-forms (*templates*) for scribe identification.

**Historical, non-cursive documents:** In contrast to word recognition methods, Alirezaee et al. developed a character recognition system for medieval Persian manuscripts. They extract statistical features from previously binarized document images. Arrivault et al. propose a combined statistical and structural character recognition approach for ancient Greek and Egyptian documents. Here, structural features such as attributed graphs are computed and classified for characters rejected during the (preceding) classification of statistical features. Another approach concerning historical Greek documents was published by Vamvakas et al. that calculates zone

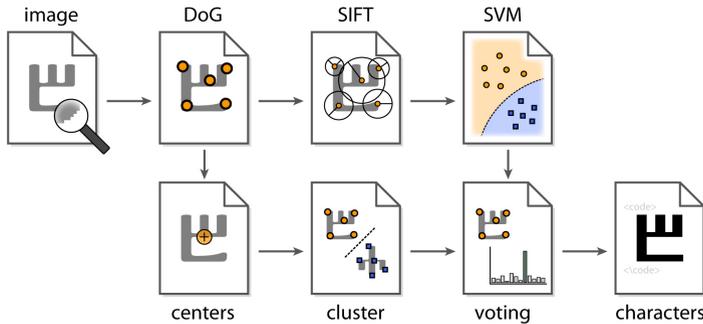


Figure 1. The proposed system consists of two task-levels: classification (upper row) and character localization (lower row).

features and character profile features on the binarized image segmented individual characters. In 2007 Ntzios et al. developed a so-called segmentation-free character recognition system applicable to the same document type. It extracts geometrical features from binarized images in combination with a watershed-like algorithm that fills cavities. A decision tree is used for the character classification.

The BIT-Alpha company (Tomasi and Tomasi 2009) presents a combined approach considering both word and character detection methods.

However, none of these methods gives positive results with faded and damaged manuscripts.

### 3. Methodology

Contrary to the methods introduced in the previous section, the proposed system here has a fundamentally new architecture, which is designed to compensate the drawbacks that arise when dealing with ancient manuscripts. Instead of applying a binarization in order to compute features, they are directly extracted from the gray-scale image.

According to its major tasks, the system is divided into two procedures: classification and localization (see Fig. 1). The fulfillment of both tasks is based upon the extraction of interest points. The interest point detector extracts blob-like regions at different scales of the manuscript image. Local descriptors robustly define the respective regions with respect to a certain set of image transformations. The descriptors are then classified by means of a multi-kernel SVM. Having classified all extracted image regions, each character consists of multiple pre-classified points. In order to assign one class label to each character present in an image, the interest points need to be clustered. *K*-means clustering groups the interest points according to the underlying characters. Finally,

a so-called interest point voting weights the class probabilities of all local descriptors belonging to the same cluster and assigns the final class label to every character.

### 3.1. Feature Extraction

As outlined, characters are detected and classified by means of interest points. These points are located at local extrema of the image's second derivative which is approximated via the Difference-of-Gaussians (DOG) function (Lowe). In other words, the interest points mark character attributes, such as junctions, endings, stroke borders, corners, and circles, and their respective size.

In order to mathematically describe the characters' attributes, local descriptors are computed at the locations of interest points. Intuitively, one could consider the gray-values of the image within a local grid (as they are the basic information which is observed by humans, too). However, this information is not robust against image transformations such as affine transformations (e.g. rotation, scale) or photometric changes (illumination, sensor noise). That is why local descriptors are computed at locations of interest points. The proposed system computes Scale-Invariant Feature Transform (SIFT) descriptors, which were first introduced by Lowe. These descriptors convert a local image grid into a 128-dimensional vector by means of gradients. Thus, the descriptors are robust against photometric changes. Additionally, SIFT descriptors are invariant with respect to rotation and scale. Considering the challenge of character recognition, it is desirable to recognize characters of different size or resolution. However, the descriptor's invariance to rotation leads to problems when recognizing characters. For instance, the Glagolitic  $\mathfrak{L}$  has the same topology as the Glagolitic  $\mathfrak{U}$ , rotated by  $180^\circ$ . If we consider descriptors that are invariant to rotation, the system cannot differentiate between these two characters. That is why the SIFT descriptors are not computed rotationally invariant, but robust against rotational changes.

Fig. 2 shows two Glagolitic characters with their corresponding interest points. Gray circles show the region of interest points which are denoted by white squares. The lines connecting the squares with the circles indicate the main orientation of each interest point. The histograms represent down-sampled local descriptors of the highlighted (black) interest point. Note that the local descriptors are the same if they are computed rotationally invariant ( $360^\circ$ ).

### 3.2. Classification

Having computed the local descriptors, character attributes – such as strokes, junctions or endings – are described by high-dimensional feature vectors. Assuming that the character attributes slightly change from one character class to another, characters can be recognized using only local information. In other words, there is no need for

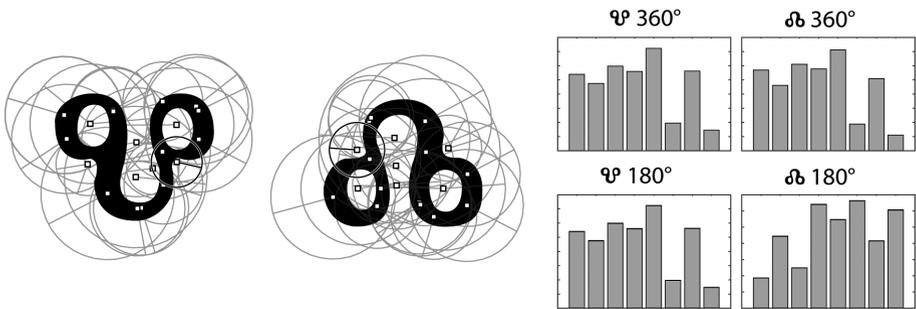


Figure 2. A Glagolitic  $\mathfrak{U}$  and  $\mathfrak{B}$  with their local descriptors (left). The down-sampled features computed rotationally invariant (right) and with rotational dependence up to 180°.

establishing a relationship between the local descriptors of one character in order to recognize the character. That is why they are directly classified in the proposed system. We use SVMs, which can be trained on classifying local descriptors – in our case character features – to assigned classes. SVMs benefit from the fact that they are based on statistical learning theory rather than error minimization. Thus, they achieve a good generalization – while still being flexible – even if a small training set is obtained. For the classification an SVM is trained using 20 manually tagged characters per character class.

In order to further improve the classification performance, one-against-all tests are performed. This means that one SVM is trained per character class. Each SVM decides whether a local descriptor belongs to its character class or not (e.g.  $\cdot\uparrow$ , not  $\cdot\downarrow$ ). In addition to the class labels predicted, a probability is assigned by each classifier resulting in a probability histogram, i.e. the assignment of the probability of each interest point of belonging to a certain character class (cf. Fig. 3). Another advantage of one-against-all tests is the fact that the classifiers are not too sensitive to noise in the training data, as the criterion function is less complex when two classes are to be considered.

### 3.3. Character Localization

For traditional OCR engines, the characters or words are localized implicitly in the binarization step. If handwriting OCR engines are considered, an additional character segmentation step needs to be performed in order to detect concatenated characters. In contrast, the proposed system has no information about the positions of characters in a given image to the point of feature classification. Indeed, the positions of the classified features are known, but – as a feature does not necessarily represent a whole character – the position and size of the character is unknown. The character localization is based

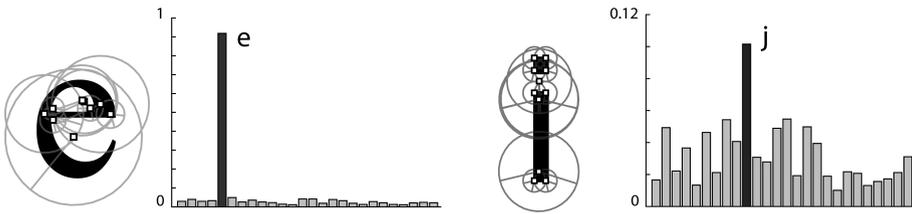


Figure 3. Probability histogram of two character clusters. A correct classification (left) and a false classification (right).

on clustering the interest points in the spatial domain, in our case applying the  $k$ -means clustering algorithm (see following paragraph). This approach benefits from the fact that degraded characters are detected with local descriptors, but not considered when the image is binarized. Thus, even degraded characters can be localized. Another advantage is the low computational complexity, since only the interest points are considered.

The  $k$ -means algorithm groups clusters of interest points and their centers. Each group should correspond to one character. But the  $k$ -means clustering cannot estimate the number of clusters  $k$ . To overcome this problem the scales of interest points are exploited. Each character produces a single local maximum in a certain scale level. When this information is extracted, the number  $k$  of the  $k$ -means can be estimated and at the same time initial cluster positions are obtained that improve convergence. Having clustered the interest points, each cluster consists of all interest points that belong to the same character.

### 3.4. Feature Voting

For the final character classification a voting scheme is applied. Therefore, all local descriptors of a cluster are considered. Each descriptor was previously classified. Hence, a probability histogram exists that indicates the class likelihood of each descriptor in the cluster. If these histograms are accumulated, the maximum bin indicates the most probable class label. Fig. 3 shows the final probability histogram of two degraded characters. Each histogram bin represents one of the previously trained character classes. The bin's height indicates each character's probability of belonging to the respective class. The left character is classified correctly, having a significantly high class probability. In contrast, the probability histogram of a false classification is given in Fig. 3 (right). There, three class probabilities are similarly high. If the histogram is indecisive, e.g. for character fragments, the alternative hypotheses could be further processed, e.g. by a dictionary.

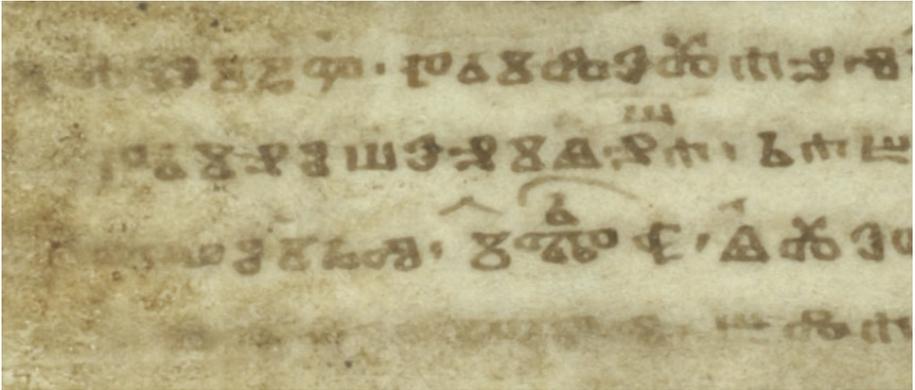


Figure 4. Random sample page of *Missale Sinaiticum*.

#### 4. Results

In this section the evaluation of the proposed system is given. In order to evaluate the system, 15 pages containing 1055 characters are extracted from the *Missale Sinaiticum*. The pages were chosen randomly (cf. Fig. 4). They contain faded-out ink, degraded characters and background noise. For groundtruthing, each character was brushed with a gray-value that corresponds to its class index.

The evaluation is based on the values of True Positives ( $TP$ ) (correctly located and correctly classified characters), False Positives ( $FP$ ) (correctly located, but falsely classified characters) and False Negatives ( $FN$ ) (characters which are not located). These values allow for computing the precision and recall. Thus, the precision indicates the percentage of correctly classified characters to those retrieved. Whereas the recall specifies the percentage of correctly classified characters to those present in an image. Mathematically, the former is defined as the sum of  $TP$  divided by the sum of retrieved values ( $TP + FP$ ). The latter is the sum of  $TP$  divided by the total number of elements that exist ( $TP + FP + FN$ ).

The aim of a classification task is to maximize both, the precision and the recall. Therefore the F score is introduced, which is a weighted average between the precision and the recall:

$$F\beta = \frac{(1 + \beta^2)p \cdot r}{\beta^2 p + r} \quad \leftrightarrow \quad F\beta = \frac{(1 + \beta^2)TP}{(1 + \beta^2)TP + \beta^2 FN + FP}$$

	F0.5-score	recall	precision	#
with clustering	0.772	0.673	0.832	1055
artificial clustering	0.804	0.748	0.837	1055

Table 1. System's recall, precision and F-score when the proposed system and artificial clustering is applied.

where  $r$  is the recall and  $p$  is the precision. The right equation expresses the F-score in terms of  $TP/FP$ . The  $\beta$  allows weighting the precision or the recall. Thus, if  $\beta$  is set to 0.5, the precision is weighted twice as much as the recall.

**System Evaluation:** In order to demonstrate the effect of the character localization, artificial clustering is implemented. This is based on the annotated ground truth where cluster centers are defined as the center-of-mass of each blob. As a constraint, only interest points being within a character blob are considered. Therefore, the character localization (clustering) does not introduce an error. Thus, the error introduced by clustering can be extracted. The system achieves an F0.5-score of 0.772 on the investigated dataset. If artificial clustering is applied, an F0.5-score of 0.805 is achieved. This directly draws the conclusion that the F-score is decreased by 0.033 because of the character localization. The test setup additionally shows that the character clustering has hardly any influence on the system's precision (difference: 0.005). In contrast, the proposed  $k$ -means decreases the recall rate by 0.075. This results from clustering errors which increase the  $FN$  rate as characters are not localized correctly.

**Evaluation of Degraded Characters:** By extracting single characters, it is possible to evaluate only the classification step illustrated in Fig. 1. Therefore two datasets are constructed that consist of single characters which were annotated and extracted from the Missal.

The first dataset (setA) consists of 10 classes having 10–12 samples (totally 107) which are well preserved. This dataset is a reference for the evaluation with degraded characters. The second dataset, which is referred to as setB, contains 25 character classes with approximately 9 characters per class (totally 198). Degraded or partially visible characters were extracted to construct this set. It is used to demonstrate the system's behavior when degraded characters need to be recognized.

Fig. 5 shows examples of both datasets. It can be seen that some characters such as  $\Omega$ ,  $\sqcup$  and  $\mathcal{U}$  are similar to each other. The degraded characters in the second row differ strongly from those of setA. They are hard to read for humans.

SetA is evaluated first in order to show the performance of the method on undistorted data. Therefore, 10 SVM kernels are trained using 10 samples per class. Then all 107 test characters are evaluated. The voting is the same as described in Section 3.4, except for

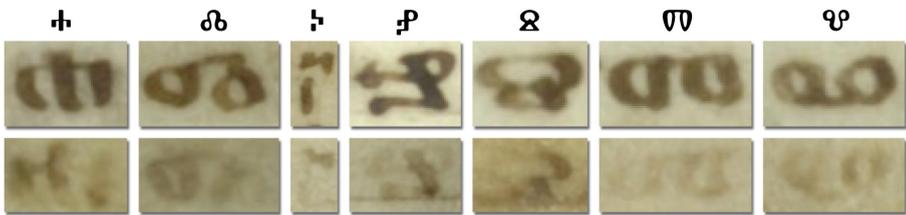


Figure 5. Examples of the datasets evaluated. The first row shows characters of setB, whereas the second row shows the same examples from the dataset containing degraded characters.

the fact that the clustering needs not to be performed. For the character classification an overall precision of 98.13% is achieved, which means that only 2 characters out of 107 are falsely predicted. Both confused characters consist of two circles and a connecting stroke (see Fig. 4 second and last column) which produce similar descriptors.

For a direct comparison of both datasets, the same ten classes were extracted from setB. Certainly the same classifier is used in both test setups. In contrast to setA the degraded characters in setB have a lower precision which is 78.89%. These numbers indicate that it is harder for the system to classify degraded characters. On the other hand the system can cope with uncertainty which arises from the fact that fewer descriptors are classified in this case.

In addition to the comparison of setA and setB, all 198 degraded characters were evaluated. Even though 25 different classes are predicted in this evaluation (+15 classes), the precision decreases slightly by 7.17%. Thus, the overall precision is 71.72% when descriptor voting is applied on degraded characters. The ratio of detected descriptors and those classified now is 26%, which means a decrease by 13% compared to the previous test on the same dataset with 10 classes. Since the performance decrease is lower than the complexity increase, the system proves to be capable for classifying degraded manuscripts.

## 5. Conclusion

This paper shows a new methodology for character recognition of ancient manuscripts. The approach, which is inspired by recent object recognition systems, exploits local descriptors directly extracted from gray-scale images. Multiple SVMs are used to classify the local descriptors. The character localization is based on clustering interest points previously extracted for the computation of local descriptors. A scale selection that adapts to the manuscript image observed allows for the cluster center initialization.

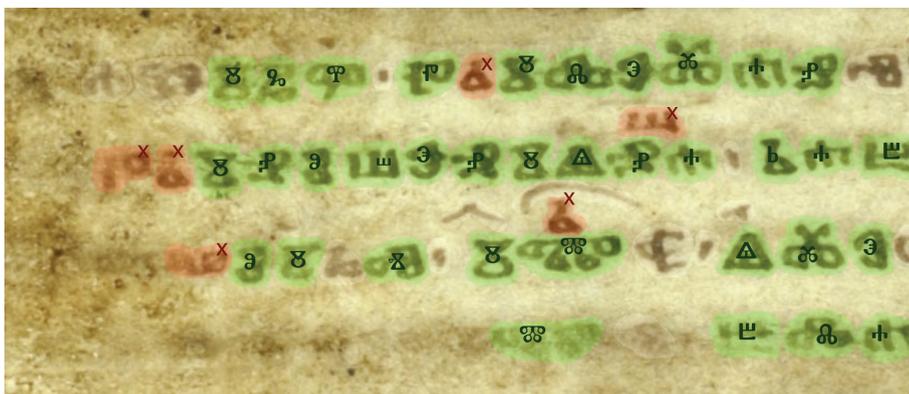


Figure 6. Sample page (cf. Fig.4) with results: Green blobs *TP*, red blobs *FP*.

The OCR system presented does not need any pre-processing of document images. In contrast to existing systems, a new architecture has been designed that focuses on images of degraded manuscripts. Since ancient manuscripts—much more often than modern—exhibit stains, faded-out ink and rippled support, they pose new challenges for OCR.

## Bibliography

- Alirezaee, Shahpour, Hassan Aghaeinia, Karim Faez and Alireza S Fard. “An Efficient Feature Extraction Method for the Middle-Age Character Recognition.” *Advances in Intelligent Computing*. Berlin/Heidelberg: Springer, 2005. 998–1006.
- Arrivault, Denis, Noël Richard, Christine Fernandez-Maloigne and Philippe Bouyer. “Collaboration between Statistical and Structural Approaches for Old Handwritten Characters Recognition.” *Graph-Based Representations in Pattern Recognition*. Eds. Luc Brun and Mario Vento. Berlin/Heidelberg: Springer, 2005. 291–300.
- Frinken, Volkmar, and Horst Bunke. “Self-Training Strategies for Handwriting Word Recognition.” *Advances in Data Mining. Applications and Theoretical Aspects*. Eds. Petra Perner. Berlin/Heidelberg: Springer, 2009. 291–300.
- Frinken, Volkmar, Tim Peter, Andreas Fischer, and Horst Bunke. “Improved Handwriting Recognition by Combining Two Forms of Hidden Markov Models and a Recurrent Neural Network.” *Computer Analysis of Images and Patterns*. Eds. André Galalowicz and Wilfried Philips. Berlin/Heidelberg: Springer, 2009. 189–196.
- Hofmeister, Wernfried, Andrea Hofmeister-Winter, and Georg Thallinger. “Forschung am Rande des paläologischen Zweifels: Die EDV-basierte Erfassung individueller Schriftzüge im Projekt DAMALS.” *KPDZ 1*. 261–92.

- Kleber, Florian, and Robert Sablatnig. “A Survey of Techniques for Document and Archaeology Artefact Reconstruction.” *10<sup>th</sup> Int. Conf. on Document Analysis and Recognition (ICDAR)*. Barcelona, 2009. (CD publication).
- KPDZ 1: *Kodikologie und Paläographie im digitalen Zeitalter – Codicology and Palaeography in the Digital Age*. Eds. Malte Rehbein, Patrick Sahle and Torsten Schaßan. Schriften des Instituts für Dokumentologie und Editorik 2. Norderstedt: Books on Demand, 2009. Online: <urn:nbn:de:hbz:38-29393>, <<http://kups.ub.uni-koeln.de/volltexte/2009/2939/>>.
- Lavrenko, Victor, Toni M. Rath, and R. Manmatha. “Holistic Word Recognition for Handwritten Historical Documents.” *1st International Workshop on Document Image Analysis for Libraries (DIAL)*, 2004. 278–287.
- Lettner, Martin, Melanie Gau, Heinz Miklas and Robert Sablatnig. “Image Acquisition & Processing Routines for Damaged Manuscripts.” *Digital Medievalist* (forthcoming).
- Lowe, David G. “Distinctive Image Features from Scale-Invariant Keypoints.” *International Journal of Computer Vision (IJCV)* 60 (2004): 91–110.
- Miklas, Heinz. “Die slavischen Schriften: Glagolica und Kyrillica.” *Der Turmbau zu Babel. Ursprung und Vielfalt von Sprache und Schrift*. Ed. Wilfried Seipel. Wien: Kunsthistorisches Museum Wien & Skira, 2003. 243–49 (No. 3.5.16–26).
- Miklas, Heinz. “Zur editorischen Vorbereitung des sog. Missale Sinaiticum (Sin. Slav. 5/N).” *Glagolitica. Zum Ursprung der slavischen Schriftkultur*. Ed. Heinz Miklas. Wien: ÖAW, 2000. 117–29, XV–XVI.
- Miklas, Heinz et al. “St. Catherine’s Monastery on Mount Sinai and the Balkan-Slavic Manuscript-Tradition.” *Slovo. Towards a Digital Library of South Slavic Manuscripts*. Sofia, 2008. 13–36, 286.
- Miklas, Heinz, Melanie Gau, Martin Lettner, and Manfred Schreiner. “Editing the Sinaitic Glagolitic Inedita. State of the Art.” *Codex Sinaiticus. Manuscripts in Modern Information Environment* (12.–13. Nov. 2009). St. Petersburg, forthcoming.
- Ntzios, Kostas, Basilios Gatos, Ioannis Pratikakis, Thomas Konidakis, and Stavros J Perantonis. “An Old Greek Handwritten Ocr System Based on an Efficient Segmentation-Free Approach.” *International Journal on Document Analysis and Recognition* 9 (2007): 179–92.
- Plamondon, Réjean, and Sargur N. Srihari. “On-Line and Off-Line Handwriting Recognition: A Comprehensive Survey.” *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22 (2000): 63–84.
- Rath, Toni M., and Rudrapatna Manmatha. “Word Spotting for Historical Documents.” *International Journal on Document Analysis and Recognition* 9 (2007): 139–52.
- Tomasi, Gilbert, and Roland Tomasi. “Approche informatique du document manuscrit.” *KPDZ 1*. 197–218.
- Vamvakas, Georgios, Basilios Gatos, Nikolaos Stamatopoulos, and Stavros J Perantonis. “A Complete Optical Character Recognition Methodology for Historical Documents.” *Document Analysis Systems (DAS)* 1 (2008): 525–32.
- Vinciarelli, Alessandro. “A Survey on Off-Line Cursive Word Recognition.” *Pattern Recognition* 35, no. 7 (2002): 1433–46.