

INTERACTIVE MACHINE LEARNING FOR WORD RECOGNITION ON  
DAMAGED HANDWRITTEN DOCUMENTS

By  
Jack Bandy

Director of Project: Brent Seales

Director of Graduate Studies: Mirosław Truszczyński

Date: April 3, 2018

MASTER'S PROJECT

Jack Bandy

The Graduate School  
University of Kentucky  
2018

INTERACTIVE MACHINE LEARNING FOR WORD RECOGNITION ON  
DAMAGED HANDWRITTEN DOCUMENTS

---

MASTER'S PROJECT

---

A document submitted in partial  
fulfillment of the requirements for  
the degree of Master of Science in  
the College of Arts and Sciences at  
the University of Kentucky

By  
Jack Bandy  
Lexington, Kentucky

Director: Dr. Brent Seales, Professor of Computer Science  
Lexington, Kentucky 2018

Copyright© Jack Bandy 2018

## ACKNOWLEDGMENTS

Acknowledge people/things here

## TABLE OF CONTENTS

Acknowledgments . . . . .	iii
Table of Contents . . . . .	iv
List of Figures . . . . .	v
List of Tables . . . . .	vi
Chapter 1 Background . . . . .	1
1.1 Related Work . . . . .	1
1.2 Motivation . . . . .	3
1.3 Contributions . . . . .	4
1.4 Literature Review . . . . .	5
Chapter 2 Methodology . . . . .	7
2.1 Data . . . . .	7
2.2 Preprocessing . . . . .	7
2.3 Labeling . . . . .	7
Chapter 3 Evaluation . . . . .	8
3.1 George Washington Dataset . . . . .	8
3.2 Wycliffe New Testament . . . . .	8
Chapter 4 The First Chapter . . . . .	9
4.1 The First Section . . . . .	9
Bibliography . . . . .	10
Vita . . . . .	13

## LIST OF FIGURES

4.1	A Simple Figure . . . . .	9
-----	---------------------------	---

## LIST OF TABLES

4.1	A Simple Table . . . . .	9
-----	--------------------------	---

## Chapter 1 Background

### 1.1 Related Work

For several decades, researchers have been developing methods for automated character and word recognition. These methods take some photograph(s) of printed or handwritten text as input, and produce a transcript of that text as output. This section provides a brief summary of methods which have influenced the course of this research area, including advances in handwriting recognition, printed text recognition, and handwritten word spotting.

The nomenclature for these related tasks can be somewhat inconsistent in the literature. For the purposes of this paper, “handwriting recognition” differs from “handwritten word spotting” in that the former aims to create full transcriptions while the latter merely locates and/or recognizes instances of a given word within a document. “Printed text recognition,” although it uses many of the same methods, refers to projects that examine machine-printed texts. The consistency of character representations and thus word representations drastically changes the task, so a distinction is necessary.

### Text Recognition

“Text recognition” here refers to recognizing *printed* texts, not handwritten texts. From a technical standpoint, automatic text recognition is the task of turning an image into the text within the image.

Object character recognition (OCR) on scans of printed documents has seen success since as early as the 1980s [1, 2]. Due to the consistency of letter shapes and sizes, fairly simple techniques for pattern recognition could accurately classify characters in the same font family. However, as early as 1987, font and size constraints were no longer needed. The authors of [3] demonstrated a system that accurately classified mixtures of dissimilar fonts of varied sizes.

Gradually, more and more constraints were eliminated. After [3] removed the need for font and size assumptions, the race was on to eliminate constraints such as alignment, color, contrast, and more. Eventually, the task of printed text recognition was one that could be done “in the wild,” [4, 5, 6] with essentially no assumptions about the nature of the text. Especially important for “in the wild” recognition was eliminating the segmentation step, as in [7], such that regions of text could be found without a processing phase devoted to localization. The ideal system, then, would be able to recognize text in any image in which a human could see text.

An important benchmark dataset for this kind of text recognition is Street View Text (SVT) [?]. SVT was harvested using pictures from Google Street View, and thus contains a heterogeneous collection of word images with a variety of fonts, colors, backgrounds, and more. (Despite the variations, word images in this set do not include handwritten characters.) The SVT dataset was released in 2010, and by



2012, [5] demonstrated state-of-the-art performance for both character recognition and word recognition by training on images from the dataset. The high degree of accuracy was achieved via unsupervised feature learning and convolutional neural networks.

Even before 2012, many researchers realized that convolutions provide an ideal mechanism for recognizing the shapes of different letters. Others have taken more general approaches to text recognition via CNNs [5, 6], . The network architectures from these papers are, on the whole, restrictively large, whereas both architectures from my experiments were able to run on my laptop.

## Handwriting Recognition

Although modern methods for printed text recognition overlap methods for handwriting recognition, especially with CNNs for “in-the-wild” recognition, the convergence happened after years of parallel research. Handwriting recognition can be divided into two major categories, “online” handwriting recognition and “offline” handwriting recognition. In the former, software tracks the location of a writing utensil as a user moves it across some surface to produce letters and words, and the precise location and motion of the utensil helps reveal the intended writing. For example, UNIPEN [8], a benchmark dataset for online handwriting recognition, includes “pen trajectory” data that specifies when and where the pen touched down and lifted up, as well as the coordinates for the path of the pen.

More relevant to this project is the task of offline handwriting recognition, in which the input comprises only a picture of the handwriting and no additional information about its creation. A canonical example of the text recognition task is the MNIST dataset [9]. MNIST comprises grayscale images of individual handwritten digits, 0 to 9, and the objective is to classify each image into the digit written inside of it. Machine learning researchers have been using this task as a benchmark for several decades [10], with error rates well below 1% since 2003 [11].

Projects using MNIST and similar datasets are premised upon many constraints. For example, a very small vocabulary or character set could be recognized if they were properly aligned and segmented, but as soon as a text ventured outside those constraints (variations on letters, misspelled words, new characters, etc.), the system would falter. Even moderately successful recognition on unconstrained datasets did not exist until the early 2000s.

This changed with the use of hidden Markov Models (HMMs) [12, 13, 14]. With statistical models built for specific languages, character and word recognition accuracies improved to over 85% (varying with respect to the test corpus). More impressively, these results came on *unconstrained* texts.

While HMMs made the way for unconstrained datasets, many demonstrations were still using the IAM dataset [15], an ad-hoc database for researchers. In other words, *truly* unrestricted handwriting recognition was still a long way off even after the strides made by HMMs. Moving forward, a collection of George Washington letters became the de-facto standard. This dataset comprised hundreds of manuscript pages from the Library of Congress, handwritten by George Washington’s secretaries.

In the mid-2000s, even state-of-the-art HMM methods yielded word error rates around 50% on datasets such as the George Washington collection. But around this time, researchers began taking a new angle at the problem. Specifically, projects focused on the process of “handwriting retrieval,” rather than attempting complete transcriptions. Such projects allow users to query a dataset of images for a given word, and essentially scans the images for visual matches of that word. For example, [16] presents a word retrieval system that achieves 63% mean average precision scores on the George Washington collection.

In [17], this approach is formalized as a viable way to generate a searchable index of handwritten papers. Their method of “wordspotting” turns the search problem into a clustering problem, where word images that are “closest” to the query word are considered matches. Wordspotting is considered more thoroughly in the following section, however it is crucial to note that this approach eliminated the need for recognizing words before retrieval. In other words, matching is done in real-time.

Building upon the success of wordspotting techniques and HMMs, [18] takes a step further and first detects *characters* in a word, before inferring a word using an ensemble of HMMs. This approach allowed the recognition of words that were never seen during training, and established new standards for the George Washington dataset.

By the time ensemble HMMs came onto the scene, neural networks were already penetrating the field of handwriting recognition [?]. By 2010, advanced techniques such as bidirectional long short-term memory (BLSTM) were successfully applied to wordspotting [?] and outperformed other methods. Finally, recurrent neural networks [19] eliminated the need for word segmentation in addition to improving state-of-the-art performance on recognition tasks.

More recently, convolutional neural networks (CNNs) have become the state-of-the-art approach for text recognition on handwritten documents [20, 21]. Many of these approaches overlap text recognition methods mentioned in the previous section, and in fact, recent neural networks are designed to recognize both printed text and handwritten text.

## Word Spotting on Damaged Handwritten Documents

In this section, the scope of related projects is narrowed down from all handwriting recognition systems, and I examine research related to word spotting on historical documents.

As previously mentioned, [17] formalized the idea of wordspotting. However, the concept was originally proposed in [?], which clustered similar words to be annotated by users, and reported success for documents written by a single person in high-quality handwriting.

### 1.2 Motivation

On the surface, optical character recognition, word recognition, and handwriting recognition appear to be solved problems. As detailed in the previous section, the

explosion of machine learning research in recent years has led to drastic improvements in performance on these tasks, and many advancements have even found their way to consumer products. For example, everyday software allows users to search within scans or photographs of printed typeface, and note-taking software can now interpret penmanship that would be indecipherable to many human readers.

However, the process of transcribing ancient documents presents a niche area of text recognition which is not addressed well by standard approaches. Many historical documents, including those reviewed in this project, were meticulously transcribed with legibility comparable to typeface, suggesting that automated transcription would be straightforward. But over time, these documents have incurred damage of all different kinds. The characters originally may have looked like typeface, but after hundreds of years of human handling, physical corrosion, chemical decay, and other processes, reading certain parts of these documents is an arduous task even for skilled textual analysts. For such cases, neither fully human transcription nor fully automated transcription is ideal.

While fully manual transcription is the most accurate solution, it is incredibly time-consuming for larger documents. Moreover, on damaged documents, skilled papyrologists are required to decipher texts. In short, human transcription is often prohibitively costly in terms of time and skilled personnel.

A fully automated transcription algorithm may successfully transcribe certain portions of a historical document, but the damaged portions can distort the algorithm’s output to the point of being unusable. This is especially true in cases where letters are literally missing. This is especially true for OCR algorithms which assume constant width, spacing, and more within a document.

An ideal solution would leverage automated transcription for the undamaged portions, and allow a human reader to fill in any gaps. I refer to this as semi-automated transcription. This project presents a pipeline for semi-automated transcription, blending the irreplicable abilities of the human eye with the efficiency and scalability of character recognition algorithms.

### **1.3 Contributions**

#### **An Interactive Approach to Automated Transcription**

In this implementation, a user first labels words or letters in the document, generating a small training set for a neural network. A trained neural network will traverse all pages of the document, recognizing occurrences of any word in its training set. If the network finds no words within an area, it documents the location as "unknown" within its output, so that a user studying the transcript can revisit the area and provide a label if possible.

Given a small set of labeled samples, train a neural network in a semi-supervised manner using both labeled and non-labeled data. Once the initial model is trained, use it to create a transcription of the full document. During the transcription process, the model keeps track of difficult word images, prioritizing them for manual labeling afterwards.

## **A Technique for Virtual Ink Restoration**

The methods used in this paper come from a variety of these related tasks, including keyword and character spotting [22, 19], word recognition [18], and handwriting recognition [23, 24].

### **1.4 Literature Review**

#### **2009**

- Finding words in alphabet soup: Inference on freeform character recognition for historical scripts [18].

#### **2012**

- A novel word spotting method based on recurrent neural networks [19].
- End-to-end text recognition with convolutional neural networks [5].

#### **2013**

- Handwritten word recognition using mlp based classifier: A holistic approach [25].
- Feature extraction with convolutional neural networks for handwritten word recognition [24].

#### **2014**

- A combined system for text line extraction and handwriting recognition in historical documents [26]

#### **2015**

- Efficient segmentation-free keyword spotting in historical document collections [7].
- Adapting off-the-shelf cnns for word spotting & recognition [22].
- Segmentation-free handwritten Chinese text recognition with LSTM-RNN [27].

#### **2016**

- On the Benefits of Convolutional Neural Network Combinations in Offline Handwriting Recognition [28].
- Reading text in the wild with convolutional neural networks [6].

- PHOCNet: A deep convolutional neural network for word spotting in handwritten documents [21].
- SpottingNet: Learning the Similarity of Word Images with Convolutional Neural Network for Word Spotting in Handwritten Historical Documents [20].

## Surveys

- A survey of document image word spotting techniques [29].
- A survey on handwritten documents word spotting [30].

## Chapter 2 Methodology

### 2.1 Data

**George Washington Dataset**

**Wycliffe Dataset**

### 2.2 Preprocessing

**Alignment**

Rotation, etc.

**Binarization**

Once the word images are segmented, the RGB image is flattened into a single grayscale channel. The image is then inverted so that the text is white and the background dark gray. To remove the gray background, a two-bin histogram is generated that heuristically determines a natural boundary between the background and foreground values. Each pixel below this boundary is set to black. FIGURE. Lastly, to deal with the varying width and height of words, each image is centered inside a canvas before being fed to the network.

**Segmentation**

To preserve the possibility of crowdsourcing the labeling task for images of individual words, each step in the pipeline assumes the document page has been segmented into words. This makes for simpler machine learning mechanisms, but adds the step of segmenting the manuscript image into individual words.

Because the Wycliffe New Testament is aligned and spaced so consistently, the segmentation for this project involved plotting a vertical projection profile of a page image to determine the location of individual lines of text. To segment the line of text into individual words, in this case, a simple threshold on the horizontal projection profile provided sufficient accuracy. Because of the alignment, the intensity vector across the line had consistent valleys at the spaces in between each word.

### 2.3 Training and Clustering

**Variational Autoencoder**

The purpose of the VAE is to learn an encoded representation of the word images without the need for ground truth label. The VAE trains by encoding its input then attempting to recreate it using the decoder. Backwards propagation of error occurs based on how similar the decoded image is to the original input image. Thus, a

successful encoding allows the decoder to accurately recreate the input image. This process is visualized in figure ??.

Figure blank showsThe experiment used two different architectures to determine the effect of expanding or contracting the size of the encoded layer, both were implemented using Keras [?] with TensorFlow [?] used as the backend.

## Clustering

To determine which word images contain the same word, clustering is performed on the encoded representation of the word images. The experiment tried two different clustering schemes: k-means and agglomerative. K-means is a popular centroid-based approach which iteratively refines the approximated center of each cluster [?, ?]. On the other hand, agglomerative clustering works by repeatedly merging the two closest clusters, resulting in a dendrogram with a single “cluster” at one end and n “clusters” at the other, where n is the number of data points.

Agglomerative clustering provides a natural fit for clustering encoded word images. Intuitively, each step groups the two most similar images into the same cluster. So instead of defining the number of clusters before the algorithm starts, the algorithm can stop when the two most similar images differ by a given amount.

## Priority Queue

In the current approach, images are labeled in the order which they appear. However, this can be modified depending on the need of the given task. For example, if a general understanding of the document’s contents is more important than a word-for-word transcription, the system could request labels for frequently occurring words. Or, for partially damaged documents, the system could request labels for “misfit” words that are particularly confusing to the encoder or to the clustering algorithm.

## Providing Labels

Labels are provided in a simple graphical user interface which displays the original word image (before inverting it and removing the background). The interface, built using TKInter, allows a user to type in a label for the word image, skip to the next word, or flag the word image as needing re-segmentation.

In the experiments, a simulated oracle was used for labeling. Because the word images were sorted in order of appearance in the text, the oracle simply labeled the image using the word at the corresponding index of the transcript.

## Transcription

A transcription of the full set of word images can be given at any time. Each word image must be encoded by the network. The encoded representation is used to determine which cluster of images the word belongs to. Once a word image is clustered, transcribing it is simply a matter of applying the label for that cluster. If the cluster is unlabeled, words in that cluster will show up in the transcript as “unknown.”

## Chapter 3 Evaluation

### 3.1 George Washington Dataset

### 3.2 Wycliffe New Testament



## Chapter 4 The First Chapter

### 4.1 The First Section

Math goes here.

Here's a figure

Figure 4.1: A Simple Figure

Here	is
a	table

Table 4.1: A Simple Table

## Bibliography

- [1] J Mantas. An overview of character recognition methodologies. *Pattern recognition*, 19(6):425–430, 1986.
- [2] VK Govindan and AP Shivaprasad. Character recognition? a review. *Pattern recognition*, 23(7):671–683, 1990.
- [3] Simon Kahan, Theo Pavlidis, and Henry S Baird. On the recognition of printed characters of any font and size. *IEEE Transactions on pattern analysis and machine intelligence*, (2):274–288, 1987.
- [4] Ray Smith. An overview of the tesseract ocr engine. In *Document Analysis and Recognition, 2007. ICDAR 2007. Ninth International Conference on*, volume 2, pages 629–633. IEEE, 2007.
- [5] Tao Wang, David J Wu, Adam Coates, and Andrew Y Ng. End-to-end text recognition with convolutional neural networks. In *Pattern Recognition (ICPR), 2012 21st International Conference on*, pages 3304–3308. IEEE, 2012.
- [6] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Reading text in the wild with convolutional neural networks. *International Journal of Computer Vision*, 116(1):1–20, 2016.
- [7] Marçal Rusiñol, David Aldavert, Ricardo Toledo, and Josep Lladós. Efficient segmentation-free keyword spotting in historical document collections. *Pattern Recognition*, 48(2):545–555, 2015.
- [8] Isabelle Guyon, Lambert Schomaker, Réjean Plamondon, Mark Liberman, and Stan Janet. Unipen project of on-line data exchange and recognizer benchmarks. In *Pattern Recognition, 1994. Vol. 2-Conference B: Computer Vision & Image Processing., Proceedings of the 12th IAPR International. Conference on*, volume 2, pages 29–33. IEEE, 1994.
- [9] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [10] Léon Bottou, Corinna Cortes, John S Denker, Harris Drucker, Isabelle Guyon, Lawrence D Jackel, Yann LeCun, Urs A Muller, Edward Sackinger, Patrice Simard, et al. Comparison of classifier methods: a case study in handwritten digit recognition. In *Pattern Recognition, 1994. Vol. 2-Conference B: Computer Vision & Image Processing., Proceedings of the 12th IAPR International. Conference on*, volume 2, pages 77–82. IEEE, 1994.
- [11] Ernst Kussul and Tatiana Baidyk. Improved method of handwritten digit recognition tested on mnist database. *Image and Vision Computing*, 22(12):971–981, 2004.

- [12] U-V Marti and Horst Bunke. Using a statistical language model to improve the performance of an hmm-based cursive handwriting recognition system. In *Hidden Markov models: applications in computer vision*, pages 65–90. World Scientific, 2001.
- [13] Horst Bunke, Samy Bengio, and Alessandro Vinciarelli. Offline recognition of unconstrained handwritten texts using hmms and statistical language models. *IEEE transactions on Pattern analysis and Machine intelligence*, 26(6):709–720, 2004.
- [14] A El-Yacoubi, Michel Gilloux, Robert Sabourin, and Ching Y. Suen. An hmm-based approach for off-line unconstrained handwritten word modeling and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(8):752–760, 1999.
- [15] U-V Marti and Horst Bunke. The iam-database: an english sentence database for offline handwriting recognition. *International Journal on Document Analysis and Recognition*, 5(1):39–46, 2002.
- [16] Toni M Rath, R Manmatha, and Victor Lavrenko. A search engine for historical manuscript images. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 369–376. ACM, 2004.
- [17] Tony M Rath and Rudrapatna Manmatha. Word spotting for historical documents. *International Journal of Document Analysis and Recognition (IJDAR)*, 9(2-4):139–152, 2007.
- [18] Nicholas R Howe, Shaolei Feng, and R Manmatha. Finding words in alphabet soup: Inference on freeform character recognition for historical scripts. *Pattern Recognition*, 42(12):3338–3347, 2009.
- [19] Volkmar Frinken, Andreas Fischer, R Manmatha, and Horst Bunke. A novel word spotting method based on recurrent neural networks. *IEEE transactions on pattern analysis and machine intelligence*, 34(2):211–224, 2012.
- [20] Zhuoyao Zhong, Weishen Pan, Lianwen Jin, Harold Mouchère, and Christian Viard-Gaudin. Spottingnet: Learning the similarity of word images with convolutional neural network for word spotting in handwritten historical documents. In *Frontiers in Handwriting Recognition (ICFHR), 2016 15th International Conference on*, pages 295–300. IEEE, 2016.
- [21] Sebastian Sudholt and Gernot A Fink. Phocnet: A deep convolutional neural network for word spotting in handwritten documents. In *Frontiers in Handwriting Recognition (ICFHR), 2016 15th International Conference on*, pages 277–282. IEEE, 2016.

- [22] Arjun Sharma et al. Adapting off-the-shelf cnns for word spotting & recognition. In *Document Analysis and Recognition (ICDAR), 2015 13th International Conference on*, pages 986–990. IEEE, 2015.
- [23] Andreas Fischer, Ching Y Suen, Volkmar Frinken, Kaspar Riesen, and Horst Bunke. A fast matching algorithm for graph-based handwriting recognition. In *International Workshop on Graph-Based Representations in Pattern Recognition*, pages 194–203. Springer, 2013.
- [24] Théodore Bluche, Hermann Ney, and Christopher Kermorvant. Feature extraction with convolutional neural networks for handwritten word recognition. In *Document Analysis and Recognition (ICDAR), 2013 12th International Conference on*, pages 285–289. IEEE, 2013.
- [25] Ankush Acharyya, Sandip Rakshit, Ram Sarkar, Subhadip Basu, and Mita Nasipuri. Handwritten word recognition using mlp based classifier: A holistic approach. *International Journal of Computer Science Issues*, 10(2):422–427, 2013.
- [26] Andreas Fischer, Micheal Baechler, Angelika Garz, Marcus Liwicki, and Rolf Ingold. A combined system for text line extraction and handwriting recognition in historical documents. In *Document Analysis Systems (DAS), 2014 11th IAPR International Workshop on*, pages 71–75. IEEE, 2014.
- [27] Ronaldo Messina and Jérôme Louradour. Segmentation-free handwritten chinese text recognition with lstm-rnn. In *Document Analysis and Recognition (ICDAR), 2015 13th International Conference on*, pages 171–175. IEEE, 2015.
- [28] Dewi Suryani, Patrick Doetsch, and Hermann Ney. On the benefits of convolutional neural network combinations in offline handwriting recognition. In *Frontiers in Handwriting Recognition (ICFHR), 2016 15th International Conference on*, pages 193–198. IEEE, 2016.
- [29] Angelos P Giotis, Giorgos Sfikas, Basilis Gatos, and Christophoros Nikou. A survey of document image word spotting techniques. *Pattern Recognition*, 68:310–332, 2017.
- [30] Rashad Ahmed, Wasfi G Al-Khatib, and Sabri Mahmoud. A survey on handwritten documents word spotting. *International Journal of Multimedia Information Retrieval*, 6(1):31–47, 2017.

## **Vita**

A brief vita goes here.