# DeepOtsu: Document Enhancement and Binarization using Iterative Deep Learning

Sheng He*, Lambert Schomaker

*ᵃBernoulli Institute for Mathematics, Computer Science and Artificial Intelligence, University of Groningen, PO Box 407, 9700 AK, Groningen, The Netherlands*

## Abstract

This paper presents a novel iterative deep learning framework and apply it for document enhancement and binarization. Unlike the traditional methods which predict the binary label of each pixel on the input image, we train the neural network to learn the degradations in document images and produce the uniform images of the degraded input images, which allows the network to refine the output iteratively. Two different iterative methods have been studied in this paper: recurrent refinement (RR) which uses the same trained neural network in each iteration for document enhancement and stacked refinement (SR) which uses a stack of different neural networks for iterative output refinement. Given the learned uniform and enhanced image, the binarization map can be easy to obtain by a global or local threshold. The experimental results on several public benchmark data sets show that our proposed methods provide a new clean version of the degraded image which is suitable for visualization and promising results of binarization using the global Otsu's threshold based on the enhanced images learned iteratively by the neural network.

*Keywords:* Document enhancement and binarization, Convolutional neural networks, Iterative deep learning, Recurrent refinement

## 1. Introduction

Extracting useful information from historical document images is a challenging problem because they usually suffer from different degradations [1], such as noise, spots, bleed-through or low-contrast ink strokes [2]. A modern retrieval system, such as the Monk system [3] which is a web-based search engine in handwritten image collections, can only provide satisfying results on high-quality handwritten images. In addition, most methods for document analysis require preprocessed and clean documents as inputs to achieve a good performance [4, 5].

---

*Corresponding author
*Email addresses:* `heshengxgd@gmail.com` (Sheng He), `L.Schomaker@ai.rug.nl` (Lambert Schomaker)
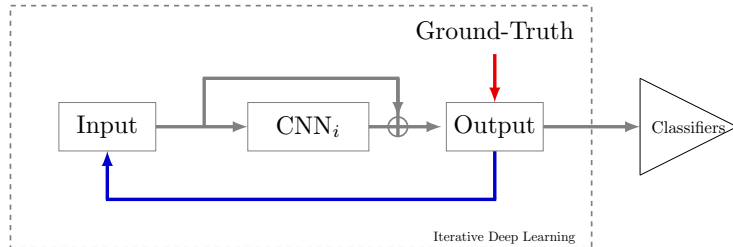
Figure 1: The proposed iterative deep learning framework. The output of the CNN$_i$ is the modified version of the input, thus it can be fed into the network iteratively for fine-tunning with different iterations. The output after several iterations which is the improved version of the input can be considered as the input of the final classifiers (SVM or Neural networks).

Document enhancement and binarization is the main pre-processing step in document analysis process. Document enhancement is the problem to improve the perceptual quality of document images and remove degradations and artifacts present in images [6], aiming at restoring its original look [7]. Document binarization is the task to separate each pixel to text and background [1]. The enhancement is also a pre-processing step for binarization on degraded document images in order to remove some unnecessary noise. Although many document enhancement methods have been proposed in the literature, most of them focus on a specific problem, such as the bleed-through correction [8, 6, 9, 10]. Few existed unsupervised methods can handle various degradations in historical documents.

Convolutional Neural Networks (CNNs) [11] has been successfully used in image classification [12] and it provides significant improvements in various applications than traditional methods. It has also been applied for document binarization [13, 14]. Since the output of binarization has the same size of the input image, the famous frameworks of neural networks are often used, such as the fully convolutional neural networks (FCNs) [15, 16], holistically-nested edge detector (HED) [17, 14] and U-Net [18]. Using deep learning provides a large margin of performance gain compared to traditional methods because millions of parameters in neural networks have been learned on a large training data set.

In this paper, we propose the iterative deep learning framework, which is shown in Fig. 1. We train the network to improve the input images, such as removing noise or correcting some degradations. Thus, the output of the neural network is the improved version of the input with supervised learning. The neural network learns the differences between the input and expected output, which might be noises or other degradations. The output can also be fed into the neural network for refinement with different iterations. After several iterations, the output which is the improved version of the input, can be used as the input for the final classifiers to improve the performance. The block of the iterative deep learning can be seamlessly integrated into any existed framework, which can be considered as the supervised data augmentation pre-processing.

In this paper, we apply the iterative deep learning for document enhancement

and the final classifier is traditional binarization method. Unlike the traditional binarization methods which train the neural network to predict the label of each pixel, we train the neural network to learn the degradation and correct the degraded document iteratively. The output of the neural network is the expected uniform and clean images instead of the binary maps, which allows the network to learn the degradations: the differences between the degraded and clean images. Since the output of the neural network is the improved version of the input, it can be also fed into the network for fine-tunning. More precisely, the learned neural network can be used recursively to refine the results because the output image can also be considered as a lightly degraded image if the learned neural network does not provide a good result in the first iteration. In addition, given the uniform image which is corrected by the learned neural network, the binarized image can be easily and efficiently obtained by a global threshold, such as the Otsu's threshold [19].

Note that the output of the proposed method is the uniform and clean version of the degraded input image, which is an acceptable view of the degraded images for the end users, such as historian, paleographers and scholars. The acceptable review is the visualization that the enhanced image should maintain the original appearance as much as possible while remove the textures and degradations on the background. Our proposed method can provide a better view of the degraded document images which only shows the original text and the noise and degradations are removed. Fig. 2 presents two examples of the original degraded images and the corresponding enhanced images of the proposed method, which shows that the enhanced images are more readable for end users than the original documents.

In summary, the differences of the proposed method with the existed works [16, 20, 14, 21] are summarized as follows. (1) Unlike the previous methods which train the neural network to learn the labels of each pixel, the output of our method is the latent uniform version of the input images, which represents an internally enhanced version of the image. (2) Our method can be used iteratively to refine the outputs since the output of the method is the improved version of the input. However, the previous methods are based on intensity probabilities per pixel, which are hard to optimize iteratively. (3) In our approach we make a distinction between handling degradations and the handling of the binarization. The neural network is trained to correct degradations, while the final binarization is achieved by the efficient global-threshold Otsu method.

The rest of this paper is organized as follows: Section 2 provides a short brief summary of a selection of related works. The proposed method is present in Section 3 and the experimental results is reported in Section 4. Finally, Section 5 provides our conclusions and prospects for future works.

## 2. Related Work

Binarization is a classical research problem for document analysis and many document binarization methods have been proposed over the past two decades in the literature. It aims to convert each pixel in a document image into either text
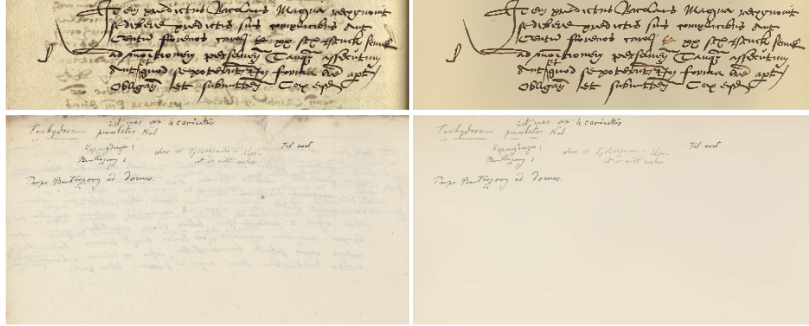
Figure 2: Demonstration of document image enhancement on two images from the Monk system [3]. The first column shows the original degraded images and the second column shows the corresponding enhanced images of the proposed method.

or background. The most popular and simple method is the Otsu [19], which is a nonparametric and unsupervised method of automatic threshold selection approach for gray-scale image binarization. It selects the global threshold based on the gray-scale histogram without any priori knowledge thus the computational complexity is linear. The Otsu method works very well on uniform and clean images while produces poor results on degraded document images with nonuniform background. In order to solve this problem, local adaptive threshold methods have been proposed, such as Sauvola [22], Niblack [23], Pai [24] and AdOtsu [25, 26]. These methods compute the local threshold for each pixel based on the local statistic information, such as the mean and standard deviation of a local area around the pixel. It should be noted that binarization is not always the goal. Methods such as Otsu can also be used for strong contrast enhancement.

Although the global or local threshold methods mentioned above are very efficient, their results are still not satisfied on high degraded and poor quality document images. Therefore, document enhancement methods are usually used as preprocessing in order to remove degradations or noise in document images. Several image processing techniques, such as the mathematical morphological operator and region-growing method are used in [27] for document enhancement and binarization. Gatos et al. [28] use a Wiener filter to estimate the background surface which is involved in the final threshold computation. Similarly, in [29], the background surface is estimated by a robust regression method and the document is binarized by a global thresholding operation. Su et al. [30] propose an adaptive contrast map for text edge detection and the local threshold is estimated based on the mean and standard deviation of pixel values on the detected edges in a local region. Nafchi et al. [31] introduce a robust phase-based binarization method which involves image denoising with phase preservation. For bleed-through correction on degraded documents, a new variational model is introduced in [6] based on wavelet shrinkage or a time-stepping scheme. A patch-based non-local restoration and reconstruction method is proposed in [32]

for degraded document enhancement. In [10], a new conditional random field (CRF)-based method [33] is proposed to remove the bleed-through from the degraded images. The bio-inspired model by the off-center ganglion cells of human vision system is used in [34] for document enhancement and binarization. All methods mentioned above use traditional techniques for document enhancement and each of them can only handle a certain type of degradation in document images.

Other priori knowledge of text is also exploit for binarization, such as the edge pixels extracted by edge detectors. For example, the Canny edge detector is used to extract edge pixels in [35] and then the closed image edges are considered as seeds to find the text region. The transition pixel which is a generation of the edge pixel is introduced in [36] is computed based on the intensity differences in a small neighbor regions and the statistic information of these pixels are used to compute the threshold. In [37], structural symmetric pixels around strokes are used to compute the local threshold. Howe [38] proposes a promising method which can tune the parameters automatically with a global energy function as a loss which incorporates edge discontinuities (Canny detector is used).

Convolutional neural networks achieve good performance on various applications, which is also applied in document analysis. For example, the winner of the recent DIBCO event [39] uses the U-Net convolutional network architecture for accurate pixel classification. In [16], the fully convolutional neural network is applied at multiple image scales. The deep encoder-decoder architecture is used for binarization in [20, 21]. A hierarchical deep supervised network is proposed in [14] for document binarization, which achieves start-of-the-art performance on several benchmark data sets. In [40], the Grid Long Short-Term Memory (Grid LSTM) network is used for binarization. However, it achieves lower performance than Vo's method [14].

## 3. Proposed Method

In this section, the problem of document enhancement is discussed based on the iterative deep learning. We first present the formulation of learning degradations for document enhancement and then the structure of CNN which is used for evaluation of the proposed model is introduced.

### 3.1. Problem formulation

An original clean or uniform document (ground-truth) is assumed to be degraded by various degradations, such as bleed-through or other artifacts. In the image enhancement formulation, the value of each pixel in the degraded images is supposed to the sum of the original value and the degraded value, which can be expressed by:

$$\mathbf{x} = \mathbf{x}_u + \mathbf{e} \tag{1}$$

where $\mathbf{x}$ is the degraded image, $\mathbf{x}_u$ is the latent clean or uniform image and $\mathbf{e}$ is the degradation. The probability density of the $\mathbf{e}$ depends on the type of degradations. Fig. 3 gives a visual example of this model.

5

Figure 3: Schematic description of the proposed degradation model. A degraded pattern $\mathbf{x}$ in the degraded image is assumed to be the sum of an ideal (uniform) pattern $\mathbf{x}_u$ and the degradation $\mathbf{e}$.

Most methods for historical document analysis requires a clean or uniform image $\mathbf{x}_u$ as input to extract the text edges or contours. Therefore, recovering the clean image $\mathbf{x}_u$ given the degraded image $\mathbf{x}$ is a classical document enhancement problem. When the uniform $\mathbf{x}_u$ is available, document binarization is quite simple, which can be computed by a global threshold, such as Otsu [19].

The Convolutional Neural Network (CNN) [11] has been successfully used in image classification [12], but it can be applied to document binarization [13, 14]. However, the traditional methods directly apply the CNN on the degraded image $\mathbf{x}$ to compute the binary image $\mathbf{x}_b$ by:

$$\mathbf{x}_b = \text{CNN}(\mathbf{x}) \tag{2}$$

which in fact requires the CNN to implicitly learn the latent uniform image $\mathbf{x}_u$, the degradation $\mathbf{e}$ and also the threshold on the latent uniform image $\mathbf{x}_u$.

In this paper, we train the neural network to predict the uniform image of the input by:

$$\mathbf{x}_u = \text{CNN}(\mathbf{x}) + \mathbf{x} \tag{3}$$

here the function $\text{CNN}(\mathbf{x}) = -(\mathbf{x} - \mathbf{x}_u) = -\mathbf{e}$ represents the degradations (negative): the difference between the degraded and clean images, which is learned by the neural network. If the input $\mathbf{x}$ is a uniform and clean image, the neural network does not need learn any information and the output of the neural network is closed to zero. Since the output $\mathbf{x}_u$ is the improved version of the input $\mathbf{x}$, it is possible to improve the output $\mathbf{x}_u$ iteratively by the neural network if we set $\mathbf{x} = \mathbf{x}_u$ in the next iteration.

When the the uniform $\mathbf{x}_u$ image is obtained after several iterations, the

binarization map can be computed by:

$$\mathbf{x}_b = \mathcal{B}(\mathbf{x}_u) \tag{4}$$

where $\mathcal{B}$ can be any existed binarization methods, or learned neural networks. Because the $\mathbf{x}_u$ is the enhanced and clean image, the simple and efficient method can be used for binarization, such as the global Otsu's threshold [19].

Our new reformulation proposed in Eqn. 3 is motivated by the residual learning [41]. However, the proposed method applies the CNN directly on the input image, which allows the neural network to learn degradations and correct the degraded images iteratively. The CNN structure in Eqn. 3 can be any neural networks, including the residual network [41].

The advantages of the proposed model are: (1) The neural network only learns the degradation of the image, without fitting the latent uniform images, such as the distribution of the background. (2) The proposed method has a new intermediate output $\mathbf{x}_u$, which can be considered as the enhanced image version of the degraded input $\mathbf{x}$. (3) This can be seamlessly integrated with other methods by Eqn. 4. For example, any existing binarization method can be applied on the estimated uniform image $\mathbf{x}_u$ learned by the neural network.

The proposed model learns the uniform image $\mathbf{x}_u$ directly from the original image. However, the learned $\mathbf{x}_u$ might not be perfect and it can also be considered as the degraded images $\mathbf{x}$ if the network does not provide good results. Thus, the learned $\mathbf{x}_u$ can be refined or enhanced recursively if we set $\mathbf{x} = \mathbf{x}_u$ to the neural network.

If we obtain the $\mathbf{x}_u$ from the neural network, there are two ways to refine it iteratively: (1) feed it into the same neural network for fine enhancement, called "Recurrent Refinement (RR)" and (2) train a new network (with the same or different structures), called " Stacked Refinement (SR)". The recurrent refinement (RR) method is defined as:

$$\mathbf{x}_u^i = \text{CNN}(\mathbf{x}_u^{i-1}) + \mathbf{x}_u^{i-1} \tag{5}$$

where $\mathbf{x}_u^i$ is the $i$-th output of the neural network and $\mathbf{x}_u^0 = \mathbf{x}$ which is the original degraded image. Note that there is only one neural network which is trained to refine the results iteratively. Fig. 4 shows the diagram of the RR framework. The advantage of the RR method is that once the neural network is trained, it can be used iteratively with many iterations. However, this also requires the neural network to learn different levels of degradations in document images. For example, it needs to remove noise on the background and recover the weak ink trace on the text region by the same network.

The stacked refinement (SR) method is defined as:

$$\mathbf{x}_u^i = \text{CNN}_i(\mathbf{x}_u^{i-1}) + \mathbf{x}_u^{i-1} \tag{6}$$

which is similar as the Eqn. 5, but the new network $\text{CNN}_i$ is trained during the $i$-th iteration to refine the input $\mathbf{x}_u^{i-1}$ which is the output of the $(i-1)$-th iteration
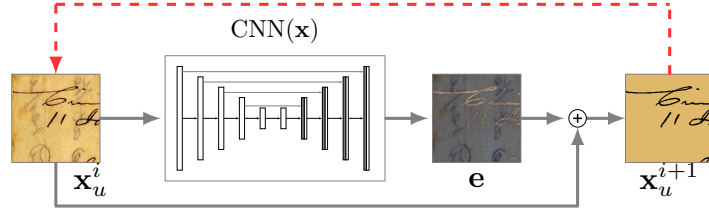
7

Figure 4: The recurrent refinement (RR) diagram of the $i$-th iteration. The red dashed line denotes that the output of the neural network can be used as input for iterative fine tuning with different iterations. $\mathbf{x}_u^0 = \mathbf{x}$ is the original degraded image at the beginning when $i = 0$.
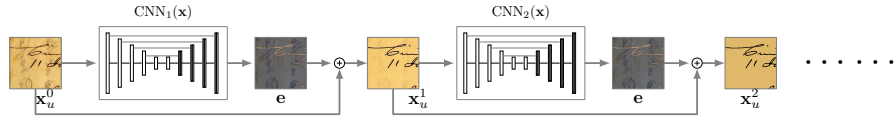


Figure 5: The stacked refinement (SR) diagram. The neural networks are stacked together to refine the results. Note that different neural networks with the same structure (CNN$_1$, CNN$_2$,...) are trained in this example.

on the CNN$_{i-1}$ neural network. Fig. 5 gives the an example of two stacked neural networks with the same structure. However, the neural network structure can also be different in different iterations. The SR method is better than the RR method because new network is trained iteratively to learn the degradations in different levels. For example, the neural network in the first iteration can learn the distribution of the background and in the second iteration can learn the distribution of the ink traces. Therefore, background noise removing and ink trace recovering can be performed in different networks.

Ideally, the RR and SR methods can be used in a mixed way. For example, the neural network of the RR method can also be a stack of networks used in the SR method which is trained iteratively. However, due to time and memory costs, we evaluate the proposed RR and SR methods separately in this paper.

### 3.2. Network Architectures

Although any neural networks can be used, in this paper we adopt the basic U-Net [18] to learn degradations in historical documents, similar to image segmentation. The architecture consists two paths: contracting path and expansive path. In the contracting path, there are five convolutional layers with the kernel size 3×3, each followed by a leaky-ReLU [42] ($\lambda = 0.25$) and 2×2 max pooling layer with stride 2. In the expansive path, the deconvoluational operation is used to upsample the feature maps and then it is concatenated with the corresponding feature maps on the contracting path, followed by a convolutional layers. The filter numbers in five convolutional layers are set to: [16,32,64,128,256], respectively. The output layer has the same size as the input image, which makes the additive operation possible (shown in Eqn. 3).

8

The network is trained by minimizing the following loss function in each iteration:

$$L^i = \frac{1}{n} \sum |\mathbf{x}_t - \hat{\mathbf{x}}_u^i| \tag{7}$$

where $n$ is the number of pixels in the image, $\mathbf{x}_t$ is the ground-truth and $\hat{\mathbf{x}}_u^i$ is the prediction from the neural network with the input $\mathbf{x}_u^i$ in the $i$-th iteration ($\mathbf{x}_u^0 = \mathbf{x}$ which is the original degraded image at the beginning). The network in each iteration is trained with the loss defined in Eqn. 7 with the degraded image and the uniform ground truth $\mathbf{x}_t$. The network of the RR and sR models on each iteration can be trained separately and jointly. In this paper, we train the network jointly and the combined loss is defined as

$$L_{total} = \frac{1}{m} \sum_i L^i \tag{8}$$

where $m$ is the number of iterations and $L^i$ is the loss on the $i$ the iteration, which is defined in Eqn. 7.

## 4. Experiments

In this section, we present the experimental results of the proposed methods for document enhancement and binarization. The training data sets are constructed based on several public benchmark data sets. We also introduce a new bleed-through data set, called Monk Cuper Set (MSC) where the historical documents are collected from the Cuper book collection of the Monk system [3].

### 4.1. Dataset

There are several public data sets for document binarization, such as (H-)DIBCO data sets which are used for document binarization competition. Similar as [14], we select images on DIBCO 2009 [43], H-DIBCO 2010 [44] and H-DIBCO 2012 [45] for training. The training set also includes documents on the Bickely-diary dataset [46], PHIDB [47] and the Synchromedia Multispectral dataset [48]. The documents on DIBCO 2011 [49],DIBCO 2013 [50], H-DIBCO 2014 [51] and H-DIBCO 2016 [52] are selected for evaluation. Four evaluation metrics which are used in the (H-)DIBCO contests, are adopted in this section for quantitatively evaluation and comparison, including F-measure, pseudo F-measure ($F_{ps}$), distance reciprocal distortion metric (DRD) and the peak signal-to-noise ratio (PSNR).

Following contest reports [43, 44, 45, 50, 51, 52], these evaluation metrics are defined as follows:

1. F-measure (FM):

$$FM = \frac{2 \times Recall \times Precision}{Recall + Precision} \tag{9}$$

where $Recall = \frac{TP}{TP+FN}, Precision = \frac{TP}{TP+FP}$, $TP, FP, FN$ denote the True positive, False position and False Negative values, respectively.

2. pseudo F-measure ($F_{ps}$):

$$F_{ps} = \frac{2 \times pRecall \times Precision}{pRecall + Precision} \quad (10)$$

where $pRecall$ is the percentage of the skeletonized ground truth image described in [44].

3. distance reciprocal distortion metric (DRD):

$$DRD = \frac{\sum_k DRD_k}{NUBN} \quad (11)$$

where $DRD_k$ is the distortion of the $k$-th flipped pixel and it is calculated using a 5×5 normalized weight matrix and $NUBN$ is the number of the non-uniform 8×8 blocks in the ground truth image (see details in [51]).

4. peak signal-to-noise ratio (PSNR):

$$PSNR = 10\log(\frac{C^2}{MSE}) \quad (12)$$

where $MSE = \frac{\sum_{x=1}^{M} \sum_{y=1}^{N} (I_{bin}(x,y) - \hat{I}_{bin}(x,y))^2}{MN}$, $C$ denotes the difference between the text and background.

We also construct a new data set for evaluation, Monk Cuper Set (MCS), which contains 25 pages sampled from a real historical collections. The documents in this set have a heavy bleed-through degradations and textural background, making them very hard for information retrieval by a computer and even hard for end users to read. Several examples are shown in Fig. 9. This data set is available on the author's website for academic usage.

*4.2. Implement details*

**Data preparation**. We train our networks with small image patches sampled from the document images with a sliding window. The basic patch size is set to 256×256 (suggested in [21]). Data augmentation is very important for boosting the performance of the neural network and we also apply augmentation methods (scale and rotation) for creating more training samples. For scale augmentation, we sample patches with the scale factor {0.75,1.25,1.5} based on the input of the neural network and resize them to 256×256. For rotation augmentation, we rotate each patch with a rotation angle 270. Overall, more than 120,000 training patches are created for training.

**Ground-truth construction**. Since the output of the neural network is the uniform images of the input, each pixel value on the ground-truth image is computed as the average pixel value with the same label within the patch. The text and background label is obtained from the ground-truth of the binary maps. For example, for the patches which do not contain any text strokes or ink traces, the ground-truth is the average image of the patch, which is helpful

Figure 6: Training samples (red box) with their corresponding ground-truth (blue box) images. Each pixel in the ground-truth is the average of pixels with the same label (text or background) within the patch.

to remove noise in the background regions of document images. Fig. 6 shows the training samples used in this paper.

**Training**. The training batch size is set to 5 due to the limitation of the memory. The learning rate is set to $10^{-4}$ and the number of training iteration is 110,000. The system runs on a PC platform with a single GPU (NVIDIA GTX 960 with 4G memory).

### 4.3. Document enhancement

Given an image, the patches with the same size as training patches are sampled with a sliding window strategy. The values of each pixel in the enhanced image are the average values of the overlapping patches computed from the trained neural networks. Fig. 7 shows a visual example of two documents (in RGB color space) on the DIBCO 2013 data set with different iterations by the SR and RR methods. It shows that with more iterations, the bleed-through and noise on background are removed and the ink traces are enhanced.

Fig. 8 and 9 show results of the enhancement documents computed by the SR method on the DIBCO 2013 and MCS data sets, respectively. From the figure we can see that: (1) The outputs are very clean. The noise and bleed-through degradations on the background are removed. (2) The large smears in the document images can not be removed completely but smoothed, because the input of the neural network is a small patch and the ground-truth of the proposed method is constructed based on the small patch, instead of the global images.

In order to quantitatively evaluate the enhancement performance, we apply two simple threshold methods to compute the binary maps: global Otsu's
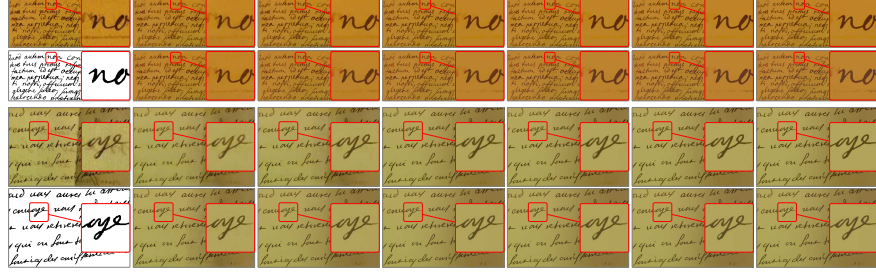
11

Figure 7: Examples of enhancement with different iterations of the SR (top row) and RR (bottom row) methods. The first column shows the original images (top) and the corresponding binary maps (bottom row). Images from second column to the last are the corresponding results of $i$-iteration where $i = 1, 2, \cdots, 6$.



Figure 8: Examples of enhancement document on the DIBCO2013 data set by the SR method.

threshold [19] and local Sauvola's threshold [22] both on the original and the enhanced images computed by the proposed methods, similar as [34]. Table 1 and 2 show the increased performance of the binarization maps computed on the original images (OI) and the enhanced images by the proposed SR and RR methods on the DIBCO 2013 and MCS datasets, respectively. From these two tables we can see that the significant improvement is achieved using the proposed deep enhancement methods as a preprocessing step. Especially, the Sauvola's threshold with the SR methods provides best performance on the two data sets, because the background of the whole image does not uniform since our model works on small patches.

Table 3 shows the performance of different methods on the MCS data set. The Sauvola's threshold based on the enhanced images produced by the SR method provides the better performance than other traditional methods on this data set. Fig. 10 presents the binary results of one document on the MCS

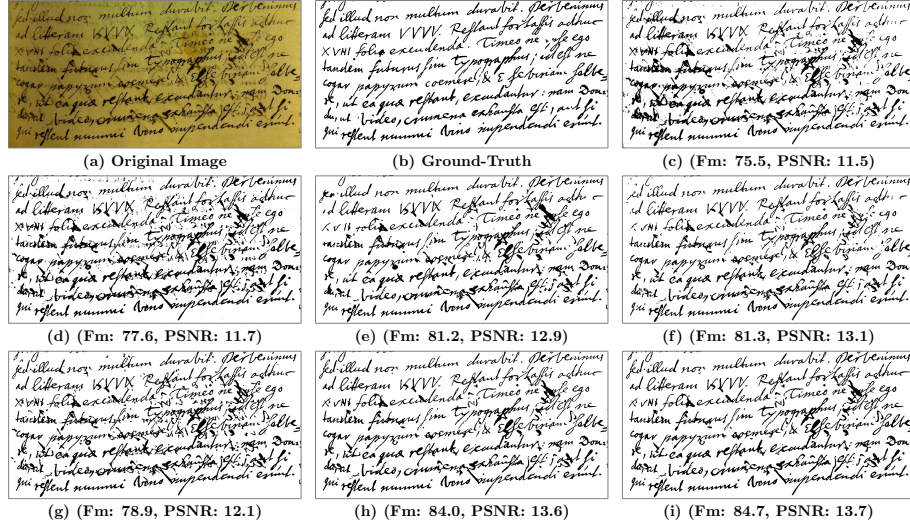Figure 9: Examples of enhancement document on the challenging MCS data set by the SR method.



|                      |                    |                             |
|----------------------|--------------------|-----------------------------|
| (a) Original Image   | (b) Ground-Truth   | (c) (Fm: 75.5, PSNR: 11.5)  |
| (d) (Fm: 77.6, PSNR: 11.7) | (e) (Fm: 81.2, PSNR: 12.9) | (f) (Fm: 81.3, PSNR: 13.1) |
| (g) (Fm: 78.9, PSNR: 12.1) | (h) (Fm: 84.0, PSNR: 13.6) | (i) (Fm: 84.7, PSNR: 13.7) |

Figure 10: Visual example with the F-measure (**Fm**) and **PSNR** metric values of the binarization results of one document on the MSC data set produced by different methods: (a) original image, (b) ground truth, (c) Otsu [19], (d) Sauvola [22], (e) Howe [38], (f) Su [53], (g) Jia [37], (h) SR-Otsu and (i) SR-Sauvola.

data set with different methods. Our proposed method provides the best visual quality and the values of evaluation metrics (F-measure and PSNR).

Table 1: Performance of binarization methods on original and enhanced images on the DIBCO 2013 data set. Values with red color show the performance improvements of the binarization results on the enhanced images (by SR and RR methods) comparing to the results on the original images (OI).

| Methods | F-measure | $F_{ps}$ | PSNR | DRD |
|---|---|---|---|---|
| OI-Otsu [19] | 80.01 | 82.82 | 16.62 | 11.00 |
| SR-Otsu | 90.00(+9.99) | 91.68(+8.86) | 20.25(+3.63) | 6.71(-4.29) |
| RR-Otsu | 88.90(+8.89) | 91.19(+8.37) | 19.62(+3.00) | 7.07(-3.39) |
| OI-Sauvola [22] | 81.23 | 83.55 | 16.60 | 11.39 |
| SR-Sauvola | 91.90(+10.67) | 93.79(+10.24) | 20.65(+4.05) | 2.60(-8.79) |
| RR-Sauvola | 90.48(+9.35) | 93.63(+10.08) | 19.97(+3.37) | 2.91(-8.48) |

Table 2: Performance of other binarization method on original and enhanced images on the MCS data set. Values with red color show the performance improvements of the binarization results on the enhanced images (by SR and RR methods) comparing to the results on the original images (OI).

| Methods | F-measure | $F_{ps}$ | PSNR | DRD |
|---|---|---|---|---|
| OI-Otsu [19] | 69.28 | 70.51 | 11.80 | 33.96 |
| SR-Otsu | 82.77(+13.49) | 85.80(+15.29) | 15.29(+3.49) | 11.32(-22.64) |
| RR-Otsu | 79.80(+10.52) | 82.31(+11.80) | 14.54(+2.74) | 15.84(-18.12) |
| OI-Sauvola [22] | 75.84 | 76.85 | 13.08 | 21.54 |
| SR-Sauvola | 87.01(+11.17) | 89.86(+13.01) | 16.19(+3.11) | 6.07(-15.47) |
| RR-Sauvola | 86.71(+10.87) | 89.68(+12.83) | 16.05(+2.97) | 6.03(-15.51) |

### 4.4. Document binarization

In this section, we first describe how to compute the binary maps given the enhanced images. Then, we present the performance of the document binarization on three benchmark datasets: DIBCO 2011, H-DIBCO 2014 and H-DIBCO 2016 data sets.

### 4.4.1. Binarization

Given the enhanced images produced by the trained neural networks, the existing method can be directly applied to compute the binary maps, such as

Table 3: Comparisons of different algorithms on the MSC data set.

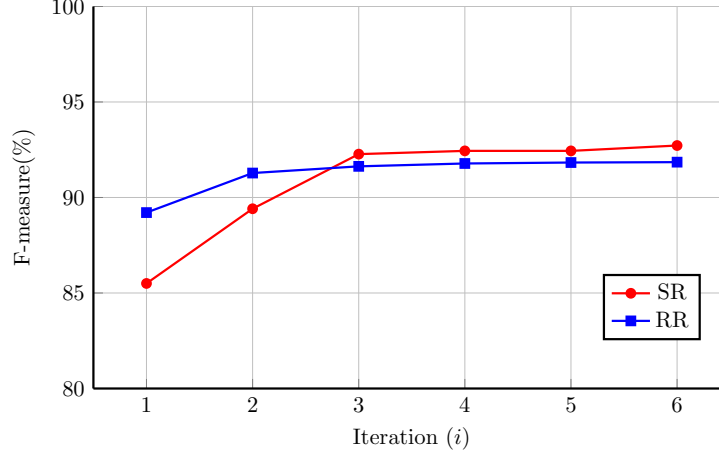| Methods | F-measure | $F_{ps}$ | PSNR | DRD |
|---|---|---|---|---|
| Otsu [19] | 69.3 | 70.5 | 11.8 | 34.0 |
| Sauvola [22] | 75.8 | 76.9 | 13.1 | 21.5 |
| Howe [38] | 85.6 | 89.1 | 15.8 | 6.4 |
| Su [53] | 82.8 | 87.4 | 15.2 | 16.8 |
| Jia [37] | 85.4 | 88.7 | 15.8 | 7.1 |
| SR-Sauvola | **87.0** | **89.9** | **16.2** | **6.1** |

Figure 11: The performance of Otsu's binarization results based on the enhanced images produced by the proposed RR and SR methods on the DIBCO 2011 data set with different iterations $i$ ($i$=1,2,...,6).

Otsu [19] which is very simple and efficient. Fig. 11 shows the performance of the Otsu's results based on the enhanced images produced by the proposed RR and SR methods with different iterations on the DIBCO 2011 data set. From the figure we can see that the trained neural network can enhance the document iteratively and the performance is dramatically increasing at the first three iterations. The performance is slightly better with more iterations.

In order to use the global Otsu's threshold, the background of the enhanced images should be uniform. However, since we use a small patch as input, the enhanced images are not global uniform if the background is nonuniform, which can be found in Fig. 8. To handle this problem, we propose several refinement steps to improve the performance of the Otsu threshold based on the enhancement images by the proposed methods.

**Uniform+Otsu**: We rescale the output of each patch to range of (0,255) if it contains ink trace, which means the background is set to values towards to 255 and the text values towards to 0. Otherwise we set it to the background. The Sauvola's threshold is used to determine whether the image patch from the output of the trained neural networks contains the ink trace or not. Fig. 12 shows an example of the Otsu results with and without using the locally uniform on the document with nonuniform background. The locally uniform is very helpful to handle the nonuniform background documents when using global Otsu's threshold.

**MS+Uniform+Otsu**: We also sample the patches with different scales on the test images (scale factors 0.75, 1.25 and 1.5 based on the size of neural network input) and resize them into 256×256. The binary map is the average of multiple scale outputs from the neural network.

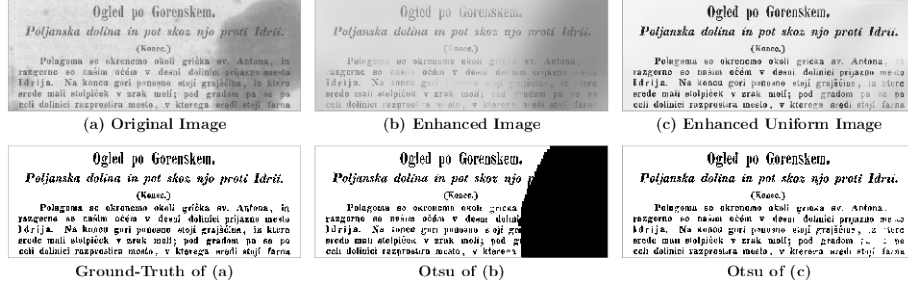**Fusion+MS+Uniform+Otsu**: Our proposed method can iteratively re-

Figure 12: Example of locally uniform results of the sample document image (PR04) on the DIBCO 2013 data set.

Table 4: Performance of binarization of the CNN output image with different refinement steps on the DIBCO 2011 data set.

| Methods | SR | | | | RR | | | |
|---|---|---|---|---|---|---|---|---|
| | F-measure | $F_{ps}$ | PSNR | DRD | F-measure | $F_{ps}$ | PSNR | DRD |
| Otsu | 92.7 | 95.6 | 19.7 | 2.2 | 91.9 | 94.9 | 19.1 | 2.6 |
| Uniform+Otsu | 92.9 | 95.7 | 19.9 | 2.1 | 92.4 | 95.3 | 19.4 | 2.4 |
| MS+Uniform+Otsu | 93.1 | 95.4 | 20.0 | 2.0 | 93.0 | 95.3 | 19.6 | 2.2 |
| Fusion+MS+Uniform+Otsu | 93.4 | 95.8 | 19.9 | 1.9 | 92.8 | 95.6 | 19.5 | 2.2 |
| Sauvola | 90.9 | 93.8 | 19.6 | 2.6 | 90.8 | 94.6 | 18.9 | 2.8 |
| Uniform+Sauvola | 93.1 | 93.7 | 19.6 | 2.2 | 92.3 | 92.6 | 19.1 | 2.7 |
| MS+Uniform+Sauvola | 92.2 | 92.2 | 19.1 | 2.5 | 91.9 | 91.7 | 18.7 | 2.7 |
| Fusion+MS+Uniform+Sauvola | 92.4 | 92.4 | 19.1 | 2.4 | 92.3 | 92.4 | 19.0 | 2.6 |

fine the outputs of the degraded inputs. However, the weak and thin ink traces might be lost at the end iteration. The performance can be improved if the outputs of each iteration are integrated. In this paper, we average the outputs of each iterations (totally six iterations) and then use Otsu to compute the binary maps.

Table 4 shows the results on the DIBCO 2011 data set. From the table we can see that using locally uniform in multiscale can improve the performance of both RR and SR methods. Combining the outputs from different iterations can provide slightly better performance for the SR method and slightly worse performance in terms of F-measure for the RR method, due to the fact the the SR method contains more convolutional layers than the RR method. We also compare the results of the Sauvola's threshold [22] based on the enhanced images computed with different refinements in Table 4. The local Sauvola's threshold provides a slightly worse performance than the global Otsu's threshold [19]. In the following section, we provide the results of both SR and RR methods using all the refinement steps (Fusion+MS+Uniform+Otsu), named as DeepOtsu.

### 4.4.2. Performance on the (H)-DIBCO benchmark data set

In this section, we provide the comparison performance of the proposed methods with other binarization algorithms on three benchmark data sets.

Table 5 shows the performance on the DIBCO 2011 data set. The Otsu's

Table 5: Comparisons of different algorithms on the DIBCO 2011 data set.

| Methods | F-measure | $F_{ps}$ | PSNR | DRD |
|---|---|---|---|---|
| Otsu [19] | 82.1 | 84.8 | 15.7 | 9.0 |
| Sauvola [22] | 82.1 | 87.7 | 15.6 | 8.5 |
| Howe [38] | 91.7 | 92.0 | 19.3 | 3.4 |
| Su [53] | 87.8 | 90.0 | 17.6 | 4.8 |
| Jia [37] | 91.9 | 95.1 | 19.0 | 2.6 |
| Vo [29] | 88.2 | 90.3 | 20.1 | 2.9 |
| Vo [14] | 93.3 | **96.4** | **20.1** | 2.0 |
| DeepOtsu(RR) | 92.8 | 95.6 | 19.5 | 2.2 |
| DeepOtsu(SR) | **93.4** | 95.8 | 19.9 | **1.9** |



Figure 13: Two failure cases on the DIBCO 2011 data set produced by the proposed DeepOtsu (SR) method. The weak strokes are missed in the binary maps.

threshold based on the enhanced images computed by the SR method provides best performance in terms of F-measure and DRD metrics, which shows that our proposed method produces the binarization maps with a less-distorted visual quality. The method proposed in [14] also uses a hierarchical neural network to predict the binary maps directly from the degraded images and it provides a slightly better performance in terms of $F_{ps}$ and PSNR. Fig. 13 provides two failure examples on the DIBCO 2011 data set of our proposed method, which misses the weak ink strokes because in the training set, the ground-truth of the weak ink strokes is also very weak which is computed by the average ink pixels in a local patch. Thus, these text regions are missed in final the binary maps computed by the Otsu's threshold. This problem could be solved by training the network with more iterations with training samples of weak ink strokes.

Table 6: Comparisons of different algorithms on the H-DIBCO 2014 data set.

| Methods | F-measure | $F_{ps}$ | PSNR | DRD |
|---|---|---|---|---|
| Otsu [19] | 91.7 | 95.7 | 18.7 | 2.7 |
| Sauvola [22] | 84.7 | 87.8 | 17.8 | 2.6 |
| Howe [38] | 96.5 | 97.4 | 22.2 | 1.1 |
| Su [53] | 94.4 | 95.9 | 20.3 | 1.9 |
| Jia [37] | 95.0 | 97.2 | 20.6 | 1.2 |
| Vo [14] | **96.7** | **97.6** | **23.2** | **0.7** |
| DeepOtsu(RR) | 94.3 | 96.3 | 20.9 | 1.9 |
| DeepOtsu(SR) | 95.9 | 97.2 | 22.1 | 0.9 |



Figure 14: Examples of the enhanced and binarization results of sample document images on the H-DIBCO 2014 data set. The left column shows the original images, the middle column shows the enhanced images by the proposed SR method and the right column shows the binarization maps based on the enhanced images.

Table 6 shows the performance of different methods on the historical document competition H-DIBCO 2014 data set. Fig. 14 shows the enhanced and corresponding binarization maps computed by the proposed DeepOtsu (SR) method on this data set. From Table 6 we can see that the performance of our proposed DeepOtsu methods is comparable to the Vo's method [14] which also uses deep learning. The performance of the traditional methods, such as Howe [38], Su [53] and Jia [37], is the comparable to the deep learning methods, such as Vo [14], which indicates that the binarization problem on the H-DIBCO 2014 data set is less challenging than other data sets. The best performance is achieved by the Vo's method [14], which integrates outputs of the binary predictions of three different networks and the final binary map is computed by a local and global threshold that is learned on a separate data set.

Table 7 presents the performance of different binarization methods on the H-DIBCO 2016 data set. Our proposed method provides better results than

Table 7: Comparisons of different algorithms on the H-DIBCO 2016 data set.

| Methods | F-measure | $F_{ps}$ | PSNR | DRD |
|---|---|---|---|---|
| Otsu [19] | 86.6 | 89.9 | 17.8 | 5.6 |
| Sauvola [22] | 84.6 | 88.4 | 17.1 | 6.3 |
| Howe [38] | 87.5 | 92.3 | 18.1 | 5.4 |
| Su [53] | 84.8 | 88.9 | 17.6 | 5.6 |
| Jia [37] | 90.5 | 93.3 | 19.3 | 3.9 |
| Vo [29] | 87.3 | 90.5 | 17.5 | 4.4 |
| Vo [14] | 90.1 | 93.6 | 19.0 | 3.5 |
| Westphal [40] | 88.8 | 92.5 | 18.4 | 3.9 |
| DeepOtsu(RR) | 90.9 | 93.9 | 19.4 | 3.1 |
| DeepOtsu(SR) | **91.4** | **94.3** | **19.6** | **2.9** |

the traditional threshold methods and also than other deep learning methods, such as the recurrent neural network [40] and the hierarchical deep supervised network [14]. Fig. 15 shows the enhanced and corresponding binarization maps computed by the proposed DeepOtsu (SR) method on this data set and we can see that the enhanced images are uniform and clean without noise and textures on the background.

From the above analysis, the proposed method works much better on document images which contains the degradation of bleed-through and noise. It provides an improved version of the degraded documents for visualization. Since the input of our method is a small patch, the large smears can not be removed completely ( Fig. 8). But this problem can be solved by training a neural network with a large input patch. The trained neural network focuses on removing the degradations, the risk is that thin or weak strokes are considered as the degradations by the neural network (Fig. 13). If reconstruction of thin strokes would be needed, this must be reflected in the composition of the training set, containing a sufficient number of such patterns.

*4.4.3. Computing time analysis*

The input patch size of the proposed is fixed, which can be computed in a very efficient way by GPU. The training takes about 24 hours for both SR and RR methods on a single GPU (NVIDIA GTX 960 with 4G memory). For testing, the computing time of the neural network for each patch is around 0.02 seconds and the processing uniform for binarization takes around 0.0008 seconds. The patches on one image can be processed in parallel on different GPUs. The training and testing time can be reduced with a stronger computer system with more GPUs and memory.

## 5. Conclusion

We have proposed a novel model for document enhancement and binarization based on iterative deep learning. Given a small patch sampled from image,
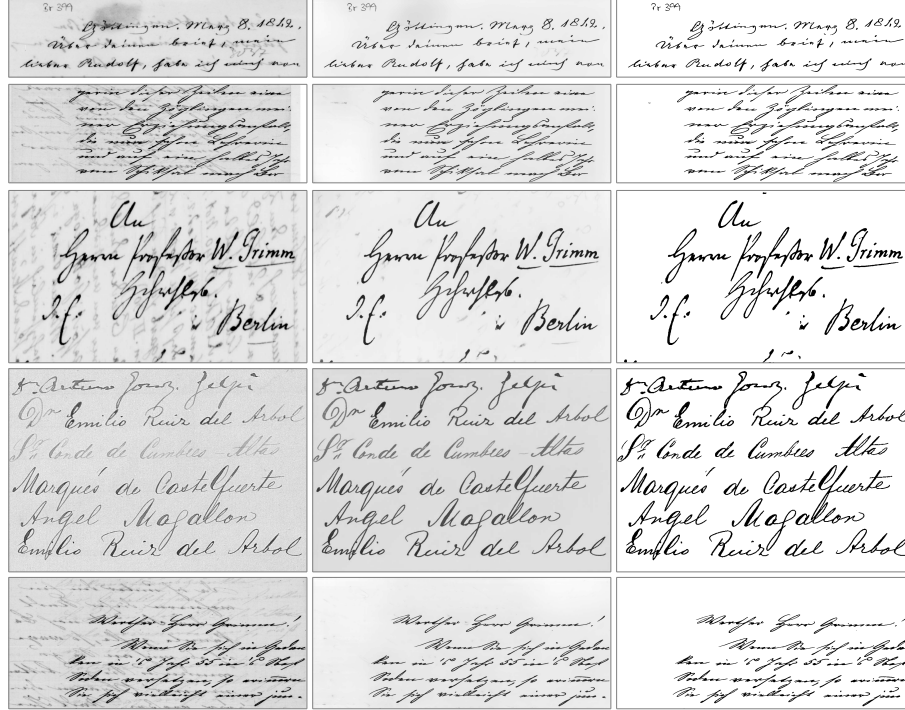
Figure 15: Examples of the enhanced and binarization results of sample document images on the H-DIBCO 2016 data set. The left column shows the original images, the middle column shows the enhanced images by the proposed SR method and the right column shows the binarization maps based on the enhanced images.

the uniform image is predicted iteratively by the proposed enhancement model in two possible ways: recurrent refinement and stacked refinement. The enhanced image produced by the proposed method is a very good view for end users, which is clean, locally uniform and does not contain any undescribable textures in the background. We evaluate the proposed method on a real historical collections from the Monk system and several public benchmark data sets. The experimental results demonstrate that our method achieves a promising performance.

In this paper, we have used the basic U-Net neural network to learn degradations in document images. More complicated neural networks can be adopted in future work, such as the ResNet [41] and DenseNet [54] and DSN used in [14]. In addition, the networks in different iterations can also be different. The complicated neural network can be applied in the beginning iterations and light neural network can be used in the end iterations for fine-tunning.

20

## Acknowledgments

## References

[1] K. Ntirogiannis, B. Gatos, I. Pratikakis, Performance evaluation methodology for historical document image binarization, IEEE Transactions on Image Processing 22 (2) (2013) 595–609.

[2] A. Tonazzini, Color space transformations for analysis and enhancement of ancient degraded manuscripts, Pattern Recognition and Image Analysis 20 (3) (2010) 404–417.

[3] T. Van der Zant, L. Schomaker, K. Haak, Handwritten-word spotting using biologically inspired features, IEEE Transactions on Pattern Analysis and Machine Intelligence 30 (11) (2008) 1945–1957.

[4] S. He, P. Samara, J. Burgers, L. Schomaker, A multiple-label guided clustering algorithm for historical document dating and localization, IEEE Transactions on Image Processing 25 (11) (2016) 5252–5265.

[5] M. Stauffer, A. Fischer, K. Riesen, Keyword spotting in historical handwritten documents based on graph matching, Pattern Recognition 81 (2018) 240–253.

[6] R. F. Moghaddam, M. Cheriet, A variational approach to degraded document enhancement, IEEE transactions on pattern analysis and machine intelligence 32 (8) (2010) 1347–1361.

[7] R. Hedjam, M. Cheriet, Historical document image restoration using multispectral imaging system, Pattern Recognition 46 (8) (2013) 2297–2312.

[8] R. F. Moghaddam, M. Cheriet, Low quality document image modeling and enhancement, International Journal of Document Analysis and Recognition (IJDAR) 11 (4) (2009) 183–201.

[9] M. R. Yagoubi, A. Serir, A. Beghdadi, A new automatic framework for document image enhancement process based on anisotropic diffusion, in: Document Analysis and Recognition (ICDAR), 2015 13th International Conference on, IEEE, 2015, pp. 1126–1130.

[10] B. Sun, S. Li, X.-P. Zhang, J. Sun, Blind bleed-through removal for scanned historical document image with conditional random fields, IEEE Transactions on Image Processing 25 (12) (2016) 5702–5712.

[11] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied to document recognition, Proceedings of the IEEE 86 (11) (1998) 2278–2324.

[12] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: Advances in neural information processing systems, 2012, pp. 1097–1105.

[13] J. Calvo-Zaragoza, G. Vigliensoni, I. Fujinaga, Pixel-wise binarization of musical documents with convolutional neural networks, in: International Conference on Machine Vision Applications (MVA), IEEE, 2017, pp. 362–365.

[14] Q. N. Vo, S. H. Kim, H. J. Yang, G. Lee, Binarization of degraded document images based on hierarchical deep supervised network, Pattern Recognition 74 (2018) 568–586.

[15] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 3431–3440.

[16] C. Tensmeyer, T. Martinez, Document image binarization with fully convolutional neural networks, in: Document Analysis and Recognition (ICDAR), 2017 14th IAPR International Conference on, Vol. 1, IEEE, 2017, pp. 99–104.

[17] S. Xie, Z. Tu, Holistically-nested edge detection, International Journal of Computer Vision 125 (1-3) (2017) 3–18.

[18] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: International Conference on Medical image computing and computer-assisted intervention, Springer, 2015, pp. 234–241.

[19] N. Otsu, A threshold selection method from gray-level histograms, IEEE transactions on systems, man, and cybernetics 9 (1) (1979) 62–66.

[20] X. Peng, H. Cao, P. Natarajan, Using convolutional encoder-decoder for document image binarization, in: Document Analysis and Recognition (ICDAR), 2017 14th IAPR International Conference on, Vol. 1, IEEE, 2017, pp. 708–713.

[21] J. Calvo-Zaragoza, A.-J. Gallego, A selectional auto-encoder approach for document image binarization, Pattern Recognition 86 (2019) 37–47.

[22] J. Sauvola, M. Pietikäinen, Adaptive document image binarization, Pattern recognition 33 (2) (2000) 225–236.

[23] W. Niblack, An introduction to digital image processing, Vol. 34, Prentice-Hall Englewood Cliffs, 1986.

[24] Y.-T. Pai, Y.-F. Chang, S.-J. Ruan, Adaptive thresholding algorithm: Efficient computation technique based on intelligent block detection for degraded document images, Pattern Recognition 43 (9) (2010) 3177–3187.

[25] R. F. Moghaddam, M. Cheriet, A multi-scale framework for adaptive binarization of degraded document images, Pattern Recognition 43 (6) (2010) 2186–2198.

[26] R. F. Moghaddam, M. Cheriet, AdOtsu: An adaptive and parameterless generalization of Otsu's method for document image binarization, Pattern Recognition 45 (6) (2012) 2419–2431.

[27] Z. Shi, S. Setlur, V. Govindaraju, Image enhancement for degraded binary document images, in: Document Analysis and Recognition (ICDAR), 2011 International Conference on, IEEE, 2011, pp. 895–899.

[28] B. Gatos, I. Pratikakis, S. J. Perantonis, Adaptive degraded document image binarization, Pattern recognition 39 (3) (2006) 317–327.

[29] G. D. Vo, C. Park, Robust regression for image binarization under heavy noise and nonuniform background, Pattern Recognition 81 (2018) 224–239.

[30] B. Su, S. Lu, C. L. Tan, Robust document image binarization technique for degraded document images, IEEE transactions on image processing 22 (4) (2013) 1408–1417.

[31] H. Z. Nafchi, R. F. Moghaddam, M. Cheriet, Phase-based binarization of ancient document images: Model and applications, IEEE transactions on image processing 23 (7) (2014) 2916–2930.

[32] R. F. Moghaddam, M. Cheriet, Beyond pixels and regions: A non-local patch means (nlpm) method for content-level restoration, enhancement, and reconstruction of degraded document images, Pattern Recognition 44 (2) (2011) 363–374.

[33] D. Song, W. Liu, T. Zhou, D. Tao, D. A. Meyer, Efficient robust conditional random fields, IEEE Transactions on Image Processing 10 (24) (2015) 3124–3136.

[34] K. Zagoris, I. Pratikakis, Bio-inspired modeling for the enhancement of historical handwritten documents, in: Document Analysis and Recognition (ICDAR), 2017 14th IAPR International Conference on, Vol. 1, IEEE, 2017, pp. 287–292.

[35] Q. Chen, Q.-s. Sun, P. A. Heng, D.-s. Xia, A double-threshold image binarization method based on edge detector, Pattern recognition 41 (4) (2008) 1254–1267.

[36] M. A. Ramírez-Ortegón, E. Tapia, L. L. Ramírez-Ramírez, R. Rojas, E. Cuevas, Transition pixel: A concept for binarization based on edge detection and gray-intensity histograms, Pattern Recognition 43 (4) (2010) 1233–1243.

[37] F. Jia, C. Shi, K. He, C. Wang, B. Xiao, Degraded document image binarization using structural symmetry of strokes, Pattern Recognition 74 (2018) 225–240.

[38] N. R. Howe, Document binarization with automatic parameter tuning, International Journal on Document Analysis and Recognition (IJDAR) 16 (3) (2013) 247–258.

[39] I. Pratikakis, K. Zagoris, G. Barlas, B. Gatos, ICDAR2017 competition on document image binarization (DIBCO 2017), in: 2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR), IEEE, 2017, pp. 1395–1403.

[40] F. Westphal, N. Lavesson, H. Grahn, Document image binarization using recurrent neural networks, in: 2018 13th IAPR International Workshop on Document Analysis Systems (DAS), IEEE, 2018, pp. 263–268.

[41] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.

[42] A. L. Maas, A. Y. Hannun, A. Y. Ng, Rectifier nonlinearities improve neural network acoustic models.

[43] B. Gatos, K. Ntirogiannis, I. Pratikakis, Icdar 2009 document image binarization contest (dibco 2009), in: Document Analysis and Recognition, 2009. ICDAR'09. 10th International Conference on, IEEE, 2009, pp. 1375–1382.

[44] I. Pratikakis, B. Gatos, K. Ntirogiannis, H-dibco 2010-handwritten document image binarization competition, in: Frontiers in Handwriting Recognition (ICFHR), 2010 International Conference on, IEEE, 2010, pp. 727–732.

[45] I. Pratikakis, B. Gatos, K. Ntirogiannis, Icfhr 2012 competition on handwritten document image binarization (h-dibco 2012), in: Frontiers in Handwriting Recognition (ICFHR), 2012 International Conference on, IEEE, 2012, pp. 817–822.

[46] F. Deng, Z. Wu, Z. Lu, M. S. Brown, Binarizationshop: a user-assisted software suite for converting old documents to black-and-white, in: Proceedings of the 10th annual joint conference on Digital libraries, ACM, 2010, pp. 255–258.

[47] H. Z. Nafchi, S. M. Ayatollahi, R. F. Moghaddam, M. Cheriet, An efficient ground truthing tool for binarization of historical manuscripts, in: Document Analysis and Recognition (ICDAR), 2013 12th International Conference on, IEEE, 2013, pp. 807–811.

[48] R. Hedjam, H. Z. Nafchi, R. F. Moghaddam, M. Kalacska, M. Cheriet, Icdar 2015 contest on multispectral text extraction (ms-tex 2015), in: Document Analysis and Recognition (ICDAR), 2015 13th International Conference on, IEEE, 2015, pp. 1181–1185.

[49] I. Pratikakis, B. Gatos, K. Ntirogiannis, Icdar 2011 document image binarization contest (dibco 2011), in: Document Analysis and Recognition (ICDAR), 2011 International Conference on, IEEE, 2011, pp. 1506–1510.

[50] I. Pratikakis, B. Gatos, K. Ntirogiannis, Icdar 2013 document image binarization contest (dibco 2013), in: Document Analysis and Recognition (ICDAR), 2013 International Conference on, IEEE, 2013.

[51] K. Ntirogiannis, B. Gatos, I. Pratikakis, Icfhr2014 competition on handwritten document image binarization (h-dibco 2014), in: Frontiers in Handwriting Recognition (ICFHR), 2014 14th International Conference on, IEEE, 2014, pp. 809–813.

[52] I. Pratikakis, K. Zagoris, G. Barlas, B. Gatos, Icfhr2016 handwritten document image binarization contest (h-dibco 2016), in: 2016 15th International Conference on Frontiers in Handwriting Recognition (ICFHR), IEEE, 2016, pp. 619–623.

[53] B. Su, S. Lu, C. L. Tan, Binarization of historical document images using the local maximum and minimum, in: Proceedings of the 9th IAPR International Workshop on Document Analysis Systems, ACM, 2010, pp. 159–166.

[54] G. Huang, Z. Liu, L. van der Maaten, K. Q. Weinberger, Densely connected convolutional networks, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, 2017, pp. 2261–2269.

[55] D. Song, D. Tao, Biologically inspired feature manifold for scene classification, IEEE Transactions on Image Processing 19 (1) (2010) 174–184.