# Light-weight Document Image Cleanup using Perceptual Loss

Soumyadeep Dey, Pratik Jawanpuria

Microsoft India

{soumyadeep.dey,pratik.jawanpuria}@microsoft.com

### Abstract

Smartphones have enabled effortless capturing and sharing of documents in digital form. The documents, however, often undergo various types of degradation due to aging, stains, or shortcoming of capturing environment such as shadow, non-uniform lighting, etc., which reduces the comprehensibility of the document images. In this work, we consider the problem of document image cleanup on embedded applications such as smartphone apps, which usually have memory, energy, and latency limitations due to the device and/or for best human user experience. We propose a light-weight encoder decoder based convolutional neural network architecture for removing the noisy elements from document images. To compensate for generalization performance with a low network capacity, we incorporate the perceptual loss for knowledge transfer from pre-trained deep CNN network in our loss function. In terms of the number of parameters and product-sum operations, our models are 65-1030 and 3-27 times, respectively, smaller than existing state-of-the-art document enhancement models. Overall, the proposed models offer a favorable resource versus accuracy trade-off and we empirically illustrate the efficacy of our approach on several real-world benchmark datasets.

## 1   Introduction

The smartphone camera have simplified the capture of various physical documents in digital form. The ease of share of digital documents (e.g., via messaging/networking apps) have made them a popular source of information dissemination. However, readability of such digitized documents is hampered when the (original) physical document is degraded. For instance, the physical document may contain extraneous elements like stains, wrinkles, ink spills, or can undergo degradation over time. As a result, while scanning such documents (e.g., via a flat-bed scanner), these elements also get incorporated into the document image. In case of capturing document images via mobile cameras, the images are prone to being impacted by shadow, non-uniform lighting, light from multiple sources, light source occlusion, etc. Such *noisy* elements not only effects the comprehensibility of the corresponding digitized document to the human readers, it may also break down the automatic (document-image) processing/understanding pipeline in various applications (e.g., OCR, bar code reading, form detection, table detection, etc). Few instances of *noisy* document images are shown in Fig. 1.

Given an input noisy document image, the aim of document image cleanup is to improve its readability and visibility by removing the noisy elements. While general (natural scene) image restoration has been traditionally explored by the computer vision community, recent works have also focused on developing cleanup techniques for document images depending on the type of noise and document-class. These include foreground background separation [23, 24, 40], differential fading problem [22, 44], removal of shadow/smear/strain [21, 22, 39, 46, 47, 49, 50], and handling ink bleed [41, 46], etc.

Recent works [6, 20, 52] view document cleanup as an image to image translation problem, modeled using deep networks. A general direction of research has been to explore deeper and more complicated networks in order to achieve better accuracy [43, 52]. However, such deep networks often require high computational resources, which is beyond many mobile and embedded applications on a computationally limited platform. Deeper networks also usually entail a higher inference time (latency), which document image processing mobile apps such as Adobe Lens, CamScanner, Microsoft Office Lens, etc., aim to minimize for best human user experience.

In this work, we propose a light-weight encoder-decoder based convolutional neural network (CNN) with skip-connections for cleaning up document images. Focusing on memory constrained mobile and embedded devices, we design a light-weight deep network architecture. It should be noted that light-weight deep network architecture usually costs generalization performance when compared with deeper networks. Hence, in order to obtain a healthy interplay between resource/latency and accuracy, we propose to employ perceptual loss function (instead of the more popular per-pixel loss function) for document image cleanup. The perceptual loss function [11] enables transfer learning by comparing high-level representation of images, obtained from a pre-trained CNNs (e.g., trained
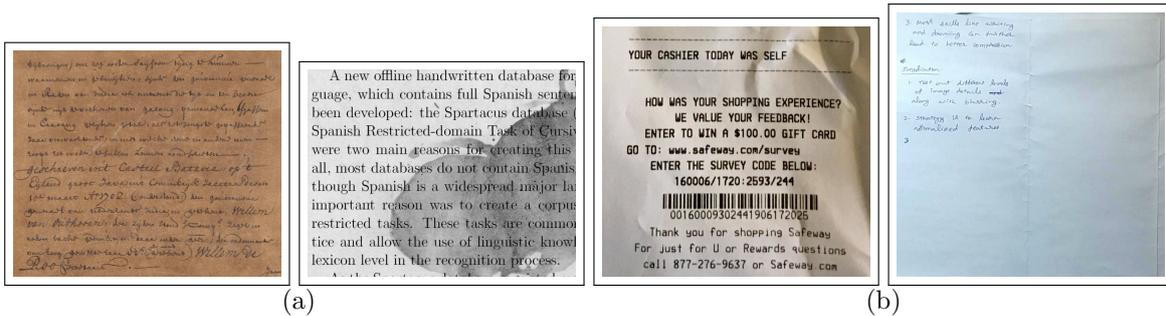
Figure 1: Typical examples of noisy document images; (a) real world degraded images from various datasets [5, 35, 51], (b) noisy real world images captured via mobile devices.

on image classification tasks). We empirically show the effectiveness of the proposed network on several real-world benchmark datasets.

The outline of the paper is as follows. We discuss the existing literature in Section 2. In Section 3, we detail our methodology. The empirical results are presented in Section 4, while Section 5 concludes the paper.

# 2    Related work

In this section, we briefly discuss existing approaches that aim to recover/enhance images of degraded documents via techniques involving binarization, and illumination/shadow correction, and deblurring, among others.

**Document image binarization**: A popular framework for document image cleanup is background foreground separation [14], where the foreground pixels are preserved and enhanced and the background is made uniform. Binarization is a technique to segment foreground from the background pixels. Analytical techniques for document image binarization involve segmenting the foreground pixels and background pixels based on some thresholding. Traditional image binarization technique such as [26] compute a global threshold assuming that the pixel intensity distribution follows a bi-modal histogram. As estimating such thresholds may difficult for degraded document images, Moghaddam and Cheriet, in [23], proposed an adaptive generalization of the Otsu's method [26] for document image binarization. In [40], Sauvola and Pietikäinen proposed a local adaptive thresholding method for the image binarization task. To improve Sauvola's algorithm's performance in low contrast setting, Lazzara and Geraud developed its multi-scale generalization in [19]. Recent works on document image binarization have also explored techniques based on conditional random fields [27], fuzzy C-means clustering [24], robust regression [48], and maximum entropy classification [22].

Deep convolutional neural networks (CNNs) have become all-pervasive in computer vision ever since AlexNet [18] won the ILSVRC 2012 ImageNet Challenge [38]. Tensmeyer and Martinez [46] posed document image binarization as a pixel classification problem and developed a fully connected convolution network for it. An encoder-decoder network was proposed in [9] to estimate the background of a document image. Then, Otsu's global thresholding technique [26] is used to obtain a binarized image with uniform background. Afzal *et al.* [1] employed a long short-term memory (LSTM) network to classify each pixel as background and foreground by considering images to be a two-dimensional sequenec of pixels. In [28], Peng *et al.* proposed a multi-resolutional attention model to learn the relationship between the text regions and background through convolutional conditional random field [16, 45]. To bypass the need of large training datasets with ground truths, Kang *et al.* [12] employed modular U-Nets [37] pre-trained for specific tasks such as dilation, erosion, histogram equalization, etc. These U-Nets are cascaded using inter-module skip connections and the final network is fine-tuned for the document image binarization task.

**Document image enhancement**: In addition to working within the background foreground separation framework, existing works have developed noise-specific document image cleanup methods such as shadow removal. Bako *et al.* [2] assumes a constant background color generates a shadow map that matches local background colors to a global reference. Similar to Bako's method, local and background colors are estimated to remove shadow from document images in [49, 50]. Inspired by the topological surface filled by water, Jung *et al.* proposed an illumination correction algorithm for document images in [39]. A document image enhancement approach have been proposed by Krigler *et al.* by representing the input image as 3D point cloud and adopting the visibility detection technique to detect the pixels to enhance [15]. Recently, Lin *et al.* [21] proposed a deep architecture to estimate (i) the global background color of the document, and (ii) an attention map which computes the probability of a pixel belonging to the shadow-free background. An illumination correction and document rectification technique using patch based encoder-decoder network is proposed in [20].

Existing works have also explored deep networks for overall document enhancement rather than focusing on correcting specific document degradations. A skip-connected based deep convolutional auto-encoder is proposed in [52]. Instead of learning the transformation function from input to output, this network learns the residual

between input and output. This residual when subtracted from the input image results in a noise free enhanced image. An end to end document enhancement framework using conditional Generative Adversarial Networks (cGAN) is proposed in [43], where an U-Net based encoder-decoder architecture is used for the generator network.

**Document image cleanup for mobile and embedded applications**: Low-resource consuming models are desirable for mobile document image processing, e.g., in apps like Adobe Lens, CamScanner, Microsoft Office Lens, etc. However, existing CNN based methods [43, 52], discussed above, propose deep architectures with huge number of parameters, making them unsuitable for memory and energy constrained devices. In this work, we propose a comparatively light-weight deep encoder-decoder based network for document image enhancement task. We employ the perceptual loss based transfer learning technique to compensate for generalization performance with a low network capacity.

# 3 Proposed approach

As discussed, we propose a light-weight deep network, suitable for mobile document image cleanup applications.

## 3.1 Network architecture

We design an encoder decoder based image to image translation network. The encoder part of the model consists of three convolution layers followed by five residual blocks. The residual blocks were first introduced in [8] for generic image processing tasks. We modify the residual blocks from the original design [8] to suite our network design. The decoder part of the model consists of five convolution layer along with skip connections from the encoder layers. This type of skip connection helps mitigate the vanishing gradient and the exploding gradient issues [8, 52, 43]. Hence, the skip connections help to simplify the overall learning of the network. Each convolution layer is followed by batch normalization layer and ReLU6 [17] activation layer.

The kernel size of the convolution layer is $3 \times 3$ and strides for all the layers is set to 1. The padding at each layer is set as *"same"*, which helps to pad the input such that it is fully covered by the filter. Padding *"same"* with stride 1 helps to keep the spatial dimension of the convolution layer output same as its input. At the end of the network a sigmoid activation function is used to obtain a normalized output between 0 and 1. The output dimension of the last layer of the decoder is either one or three depending on the end task of the network. If the network is trained for the task of binarization or gray scale cleanup then the output dimension of the last layer is set to one. For color cleanup task the output dimension of the last layer is set to three.

We term our models as M-x, where x represents the value of the maximum width of the network. In our experiments, we have considered x to be 16, 32 and 64. The M-64 model is shown in Fig. 2. In M-32 model's architecture, the output dimensions of the residual blocks are 32 and the CNN blocks with output dimensions 64 are removed. Similarly, in the case of M-16, the output dimensions of the residual blocks are set to 16 and the CNN blocks with output dimensions 32 and 64 are removed.

## 3.2 Loss function

The network is optimized by minimizing the loss function $L$ computed using Eq. 1:

$$L(I_t, I_g) = \lambda_1 \ell_1(I_t, I_g) + \lambda_2 \ell_2(I_t, I_g) + \lambda_3 \ell_3(I_t, I_g). \tag{1}$$

Here, $\ell_1(I_t, I_g) = \|I_t - I_g\|_1$ is the 1-norm loss between the translated image $I_t$ and the ground truth image $I_g$ in $YC_bC_r$ color space for color image cleanup. For gray scale cleanup, $\ell_1(I_t, I_g) = \|I_t - I_g\|_1$, refer to the 1-norm loss between the two images in gray scale. In addition to pixel-level loss function $\ell_1(I_t, I_g)$, we also employ the perceptual loss functions $\ell_2(I_t, I_g)$ and $\ell_3(I_t, I_g)$ in Eq. 1.

Perceptual loss functions [11, 36, 53] compute the difference between images $I_t$ and $I_g$ at high level feature representations extracted from a pre-trained CNNs such as those trained on ImageNet image classification task. They are more robust in computing distance between images than pixel-level loss functions. In the context of developing light-weight document image cleanup models, perceptual loss functions serve an additional role of enabling transfer learning. The perceptual loss functions in Eq. 1 helps to transfer the semantic knowledge already learned by the pre-trained CNN network to our smaller network.

The perceptual loss has two components [11]: feature reconstruction loss $\ell_2(I_t, I_g)$ and style loss $\ell_3(I_t, I_g)$. Feature reconstruction loss encourages the transformed image to be similar to ground truth image at high level feature representation as computed by a pre-trained network $\wp$. Let $\wp_j(I)$ be the activations of the $j^{th}$ layer of the pre-trained network $\wp$. Then, Eq 2 represents the feature reconstruction loss:

$$\ell_2(I_t, I_g) = \frac{1}{H_j W_j C_j} \|\wp_j(I_t) - \wp_j(I_g)\|_1, \tag{2}$$

where the shape of $\wp_j(I)$ is $H_j \times W_j \times C_j$. The feature reconstruction loss penalizes the transformed image when it deviate from the content of the ground truth image. Additionally, we should also penalize the transformed
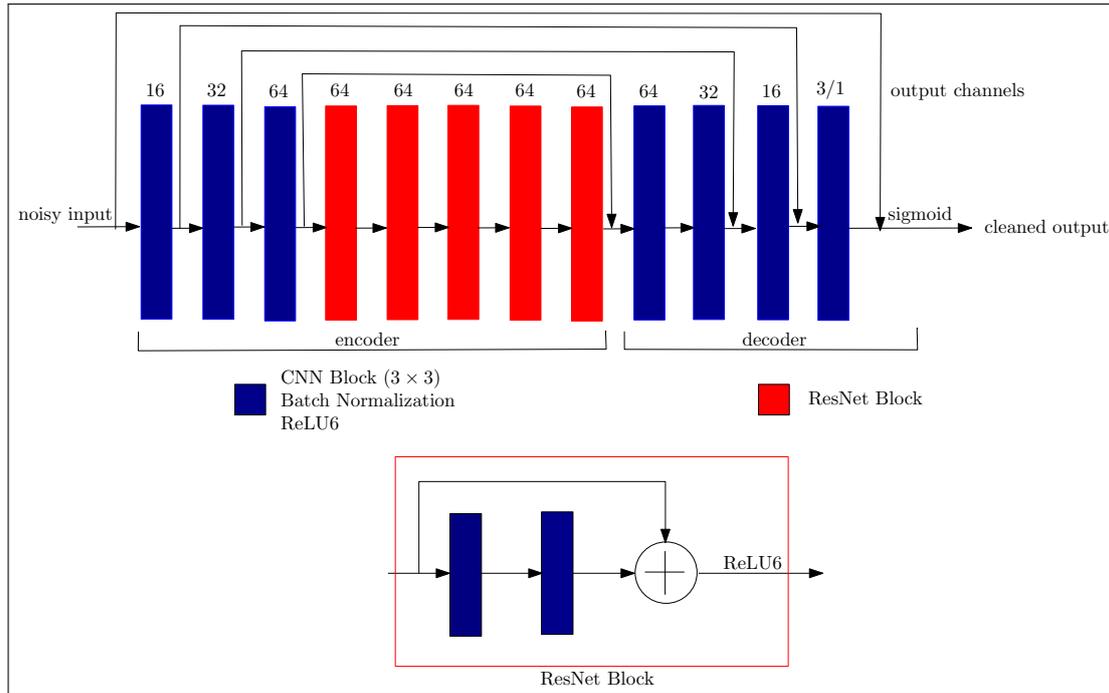
Figure 2: Proposed light-weight CNN architecture used for document image cleanup

Table 1: No. of parameters, product-sum operations, and other statistics of various models. Average per-patch inference time is reported. As an example, a $2560 \times 2560$ image has 100 patches.

| Model | Mult-Adds (in billions) | Parameters (in millions) | Size (in KB) | Inference time (in seconds) | Load time (in seconds) |
|---|---|---|---|---|---|
| DE-GAN [43] | 46.1 | 31.00 | 121 215 | 0.36 | 9.96 |
| SkipNetModel [52] | 106.8 | 1.64 | 6380 | 0.34 | 3.62 |
| M-64 (proposed) | 15.1 | 0.46 | 1779 | 0.24 | 3.25 |
| M-32 (proposed) | 6.7 | 0.11 | 445 | 0.07 | 2.89 |
| M-16 (proposed) | 1.7 | 0.03 | 111 | 0.02 | 2.68 |

image if it deviate from the ground truth image in terms of common feature, texture, etc. To achieve this style loss is incorporated as proposed in [11]. The style loss is represented in Eq. 3 as follows:

$$\ell_3(I_t, I_g) = \sum_{\forall j \in J} \|G_j^{\wp}(I_t) - G_j^{\wp}(I_t)\|, \tag{3}$$

where $\wp$ represent pre-trained CNN network, $J$ represent set of layers of $\wp$ used to compute style loss, and $G_j^{\wp}(I)$ represent a Gram matrix containing second-order feature covariances. Let $\wp_j(I)$ be the activation of the $j^{th}$ layer of the pre-trained network $\wp$, where the shape of $\wp_j(I)$ is $H_j \times W_j \times C_j$. Then, the shape of the Gram matrix $G_j^{\wp}(I)$ is $C_j \times C_j$ and each element of $G_j^{\wp}(I)$ is computed according to Eq 4 as follows:

$$G_j^{\wp}(I)_{c,c'} = \frac{1}{H_j W_j C_j} \sum_{h=1}^{H_j} \sum_{w=1}^{W_j} \wp_j(I)_{h,w,c} \wp_j(I)_{h,w,c'}. \tag{4}$$

In our work, we use the $VGG19$ network [42] trained on the ImageNet classification task [4] as our pre-trained network $\wp$ for the perceptual loss. Here, feature reconstruction loss is computed at layer conv1-2 and style reconstruction loss is computed at layers conv1-1, conv2-1, conv3-1, conv4-1, and conv5-1.

# 4  Experimental results and discussion

We evaluate the generalization performance of the proposed models on binarization, gray scale, and color cleanup tasks.

**Experimental setup**: In our experiment, the input of the network is set as $256 \times 256$. The input to the network is a 3 channel RGB image, whereas the output dimension of the network is set as 1 or 3 depending on the downstream

Table 2: Results on DIBCO13 [32]

| Model | F-measure | $F_{ps}$ | PSNR | DRD |
|---|---|---|---|---|
| Otsu [26] | 83.9 | 86.5 | 16.6 | 11.0 |
| Sauvola *et.al.* [40] | 85.0 | 89.8 | 16.9 | 7.6 |
| Tensmeyer *et.al.* [46] | 93.1 | 96.8 | 20.7 | 2.2 |
| Vo *et.al.* [48] | 94.4 | 96.0 | 21.4 | 1.8 |
| DE-GAN [43] | 99.5 | 99.7 | 24.9 | 1.1 |
| SkipNetModel [52] | 95.3 | 96.6 | 22.8 | 1.5 |
| M-64 (proposed) | 94.1 | 95.7 | 21.7 | 2.1 |
| M-32 (proposed) | 92.3 | 93.3 | 20.4 | 2.5 |
| M-16 (proposed) | 90.4 | 91.6 | 19.9 | 3.1 |

task. If the downstream task is to obtain an image in gray scale or a binary image, then the output dimension is set as 1. The output dimension is set as 3 for color cleanup task.

To handle different type of noise at various resolution, the training images are scaled at scale 0.7, 1.0, and 1.4. Further at each scale, the training images are divided into overlapping blocks of $256 \times 256$. During training, a few random patches from the training images are also used for data augmentation using random brightness-contrast, jpeg noise, ISO noise, and various types of blur [3]. Randomly selected 80% of the training patches is used to train the model while the remaining 20% is kept for validation. The model with best validation performance is saved as the final model. The network is optimized using Adam algorithm [13] with default parameter settings. The parameters $\lambda_1$, $\lambda_2$, and $\lambda_3$ of Eq 1 is set to $1e1$, $1e-1$, and $1e1$ respectively. During inference, an input image is divided into overlapping blocks of $256 \times 256$. Each patch is inferred using the trained model. Finally, all the patches are merged to obtain the final result. We use simple averaging for the overlapping pixels of the patches.

**Compared algorithms**: We compare the proposed models with recently proposed deep CNN based document image cleanup models: SkipNetModel [52] and DE-GAN [43]. Table 1 presents a comparative analysis of our proposed models with SkipNetModel and DE-GAN in terms of: (i) number of multiplication and addition operations (Mult-Adds) associated with the model [10], (ii) number of parameters, (iii) actual size on device, (iv) model load time, and (v) model inference time. A comparison with respect to these parameters is essential if the applicability of any model for memory and energy constrained devices is to determined. We implemented the models using TensorFlow Lite (https://www.tensorflow.org/lite) on an Android device with Qualcomm SM8150 Snapdragon 855 chipset and 6GB RAM size. On the device, we observe that our models are 65-1090 and 3-55 times lighter in size than DE-GAN and SkipNetModel, respectively. Similarly, our models has lesser product-sum operations and prediction time during the inference stage, making them suitable to mobile and embedded applications.

## 4.1 Binarization

We begin by discussing our results on document image binarization task. For our experiment, we have considered the publicly available binarization dataset DIBCO13 [32] and DIBCO17 [35] as test sets. The proposed models and SkipNetModel are trained on the datasets [7], [29], [30], [31], [25], [34] and [33]. While training the models for the test set DIBCO13 [32], we also include the dataset DIBCO17 [35] into our training data. The models for this task are trained using the augmentation strategy described in Sec. 4. The same training strategy is also followed while training the models for the task DIBCO17 [35]. The models are compared using the DIBCO13 [32] evaluation criteria: F-measure, pseudo F-measure ($F_{ps}$), peak signal to noise ratio (PSNR), and distance reciprocal distortion (DRD). For the metrics F-measue, $F_{ps}$, and PSNR, higher values correspond to better performance whereas, in
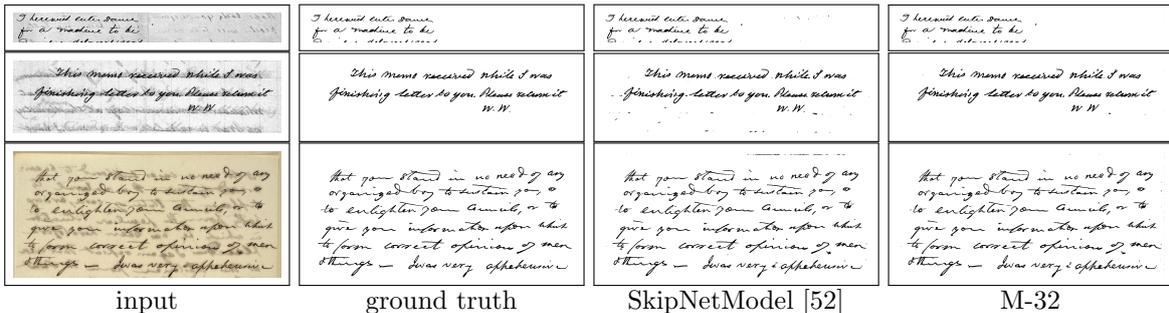


| input | ground truth | SkipNetModel [52] | M-32 |

Figure 3: Typical examples of DIBCO13 [32].

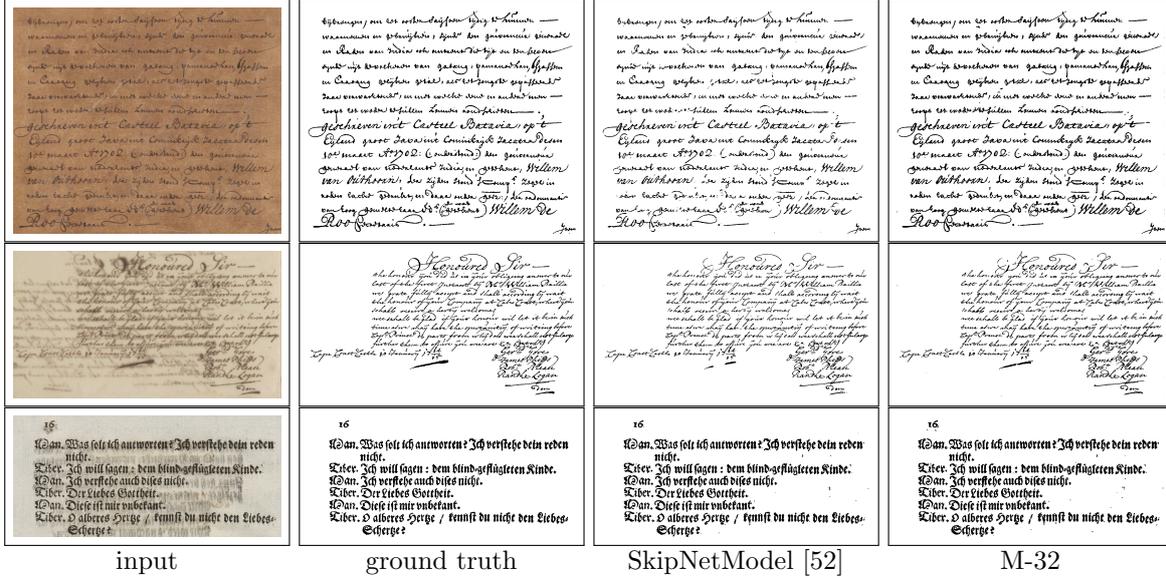| input | ground truth | SkipNetModel [52] | M-32 |

Figure 4: Examples from DIBCO17 [35].

Table 3: Results on DIBCO17 [35]. Here 10, 17a, 12, 1b, and 1a are the top 5 methods from DIBCO 2017 competition [35]

| Model | F-measure | $F_{ps}$ | PSNR | DRD |
|---|---|---|---|---|
| 10 [35] | 91.04 | 92.86 | 18.28 | 3.40 |
| 17a [35] | 89.67 | 91.03 | 17.58 | 4.35 |
| 12 [35] | 89.42 | 91.52 | 17.61 | 3.56 |
| 1b [35] | 86.05 | 90.25 | 17.53 | 4.52 |
| 1a [35] | 83.76 | 90.35 | 17.07 | 4.33 |
| DE-GAN [43] | 97.91 | 98.23 | 18.74 | 3.01 |
| SkipNetModel [52] | 91.13 | 92.91 | 18.01 | 3.22 |
| M-64 (proposed) | 90.80 | 91.73 | 17.84 | 3.32 |
| M-32 (proposed) | 89.93 | 90.61 | 17.32 | 3.74 |
| M-16 (proposed) | 87.81 | 89.40 | 16.91 | 4.15 |

case of the metric DRD lower is better. While evaluating our methods on DIBCO13 dataset, we have compared our methods with traditional binarization algorithms [26, 40], state of the art binarization techniques [46, 48], DE-GAN [43] and SkipNetModel [52]. Overall performance of these methods are reported in Table 2. In this table, performance of the methods [26, 40, 43, 46, 48] are reported as they are reported in [43]. From this table, it is evident that DE-GAN outperforms all other methods in terms of all metrics. However, the proposed method performs better than the traditional binarization algorithms [26, 40] and they perform more or less similar to other state of the art techniques. Moreover, from Tables 1 and 2, we can observe that though the proposed method can not outperform the state of the art techniques but they perform similar to most of the state of the art techniques with much lesser computational and memory cost. We have also reported the performance of the proposed methods with the top 5 methods of DIBCO17 competition [35], SkipNetModel and DE-GAN in Table 3. A similar performance of the proposed methods is also observed from this table in comparison to the state of the art techniques. Typical examples from the datasets DIBCO13 and DIBCO17 are shown in Figs. 3 and 4.

## 4.2 Gray scale and color cleanup

For the purpose of document cleanup, we first show the effectiveness of the proposed methods in gray scale. **The gray scale cleanup** part of our experiment is conducted on the publicly available dataset *NoisyOffice* [5, 51] . This dataset consists of two parts, first one is real noisy images consisting of 72 files, and a synthetic dataset consisting of 216 files. There were no groundtruth images available for the real data, therefore, we are not able to include the real dataset for quantitative analysis of our experiment. The model is trained and evaluated on the synthetic data. We divide the synthetic data into two parts - 172 images for training and 44 images for testing. The authors of [43] did not share the saved model for gray scale cleanup. Therefore, for this experiment, only the method proposed in [52] is used for comparison. To measure the capability of the proposed models with respect
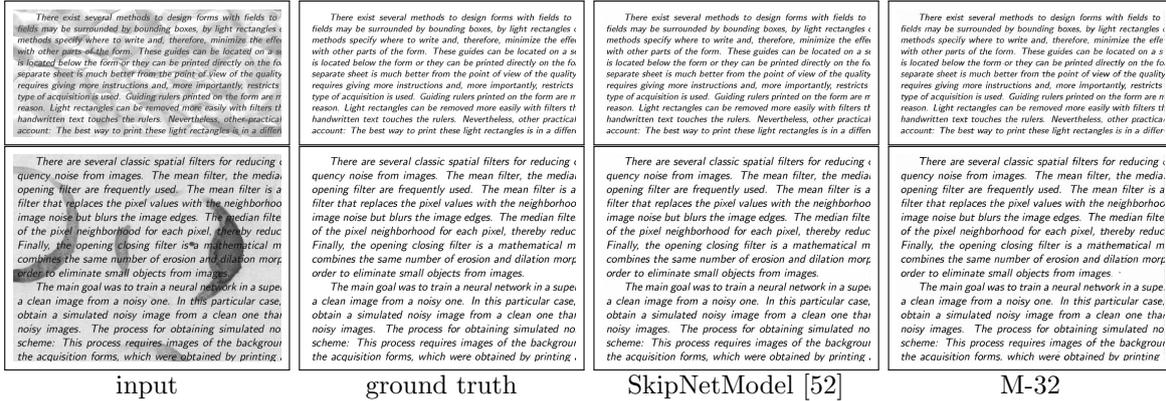
Figure 5: Typical examples of noisy images from our test set of synthetic data from [51].
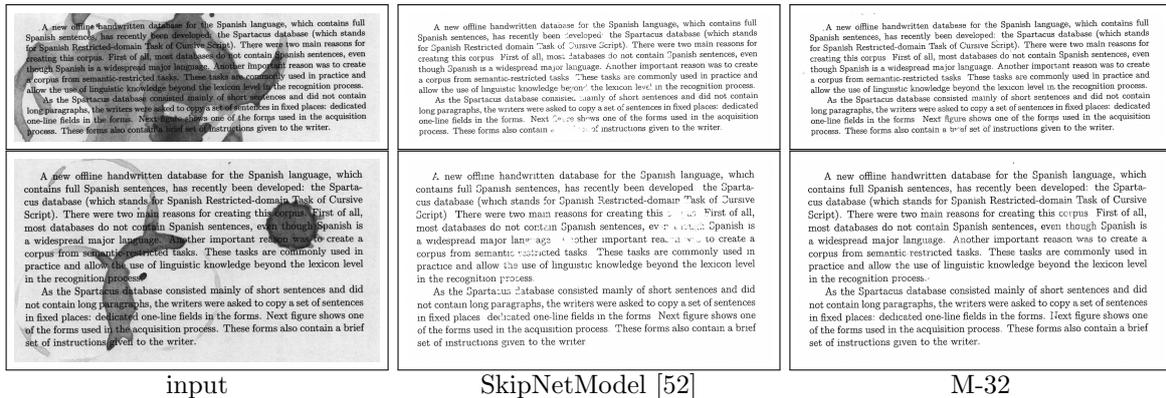


Figure 6: Typical examples of noisy images of real data from [51].

to removing noise, we adopt peak signal to noise ratio (PSNR) as the quality metric. In order to determine the dependence of the model performance on the amount of available training data, we have trained the SkipNetModel and the proposed models M-64, M-32, and M-16 by varying the amount of training data from 25% to 100%. The performance of the models with respect to PSNR score is shown in Fig. 7. It can be observed from this figure that in presence of 100% training data, SkipNetModel [52] performs better than the proposed models. However, it can also be observed from this figure that the performance of the proposed models is more or less remains the same. It can also be observed from this figure that the performance of SkipNetModel varies a lot with the variation in the amount of training data. Typical examples of inputs, groundtruths from the test set along with the outputs of the models trained on 100% training data are shown in Fig. 5. We have also shown a few examples of inputs and outputs of the trained models on real data in Fig. 6. It can be seen from this figure that the model M-32 performs better than SkipNetModel in few of the examples. From Figs. 6 and 7, we can conclude that the proposed model is more generalized and performs more robustly with respect to SkipNetModel.

Finally, we present our experimental results with respect to document **color cleanup** task. One challenging aspect of color cleanup is the preservation of color of the foreground pixels. For this experiment, we used a color dataset consisting of 250 mobile captured images. Each of the images are manually cleaned. We followed the same training strategy for training our model as described in Sec. 4. Random real life images (not belonging to the train/test set) and their corresponding outputs are shown in Fig. 8. From this figure, we can observe a decent performance of the proposed model in performing color cleanup of document images. However, to provide a quantitative measure of our method, we compute PSNR, and structural similarity index (SSIM) score on the test

Table 4: Color cleanup performance of the proposed model: SSIM and PSNR score of the noisy input images with respect to the groundtruth are 0.87 and 16.3 respectively.

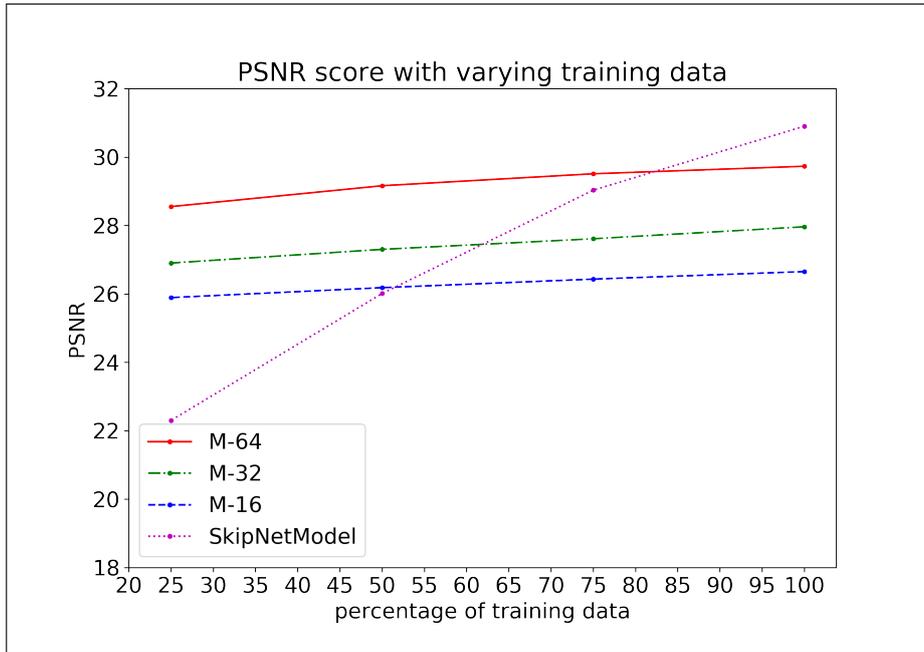| Model | SSIM | PSNR |
|---|---|---|
| M-64 (proposed) | 0.967 | 22.8 |
| M-32 (proposed) | 0.950 | 21.4 |
| M-16 (proposed) | 0.923 | 19.8 |

Figure 7: PSNR scores of the models on NoisyOffice dataset [51] with varying the training data from 25% to 100%

set of the data in Table 4.

## 5   Conclusion

We have proposed an encoder-decoder based document cleanup model for resource constrained environments. To this end, we design a light-weight deep network with only a few residual blocks and skip connections. Our loss function incorporates the perceptual loss, which enables transfer learning from pre-trained deep CNN networks.

We develop three models based on our network design, with varying network width. In terms of the number of parameters and product-sum operations, our models are 65-1030 and 3-27 times, respectively, smaller than a recently proposed GAN based document enhancement model [43]. In spite of our relatively low network capacity, the generalization performance of our models on various benchmarks are encouraging and comparable with several document image cleanup techniques with deep architectures such as [46, 52]. In addition, our models are more robust to low training data regime than [52]. Hence, the proposed models offer a favorable trade-off between memory/latency and accuracy, making them suitable for mobile document image cleanup applications.

## References

[1] Afzal, M.Z., Pastor-Pellicer, J., Shafait, F., Breuel, T.M., Dengel, A., Liwicki, M.: Document image binarization using lstm: A sequence learning approach. In: Proc. of the 3rd Int. Workshop on Historical Document Imaging and Processing. p. 79–84. HIP '15, Association for Computing Machinery, New York, NY, USA (2015)

[2] Bako, S., Darabi, S., Shechtman, E., Wang, J., Sunkavalli, K., Sen, P.: Removing shadows from images of documents. Asian Conf. on Computer Vision (ACCV 2016) (2016)

[3] Buslaev, A., Iglovikov, V.I., Khvedchenya, E., Parinov, A., Druzhinin, M., Kalinin, A.A.: Albumentations: Fast and flexible image augmentations. Information **11**(2) (2020)

[4] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database. In: CVPR09 (2009)

[5] Dua, D., Graff, C.: UCI machine learning repository (2017)

[6] Gangeh, M.J., Tiyyagura, S.R., Dasaratha, S.V., Motahari, H., Duffy, N.P.: Document enhancement system using auto-encoders. In: Workshop on Document Intelligence at NeurIPS 2019 (2019)

[7] Gatos, B., Ntirogiannis, K., Pratikakis, I.: Icdar 2009 document image binarization contest (dibco 2009). In: 10th Int. Conf. on Document Analysis and Recognition. pp. 1375–1382 (2009)
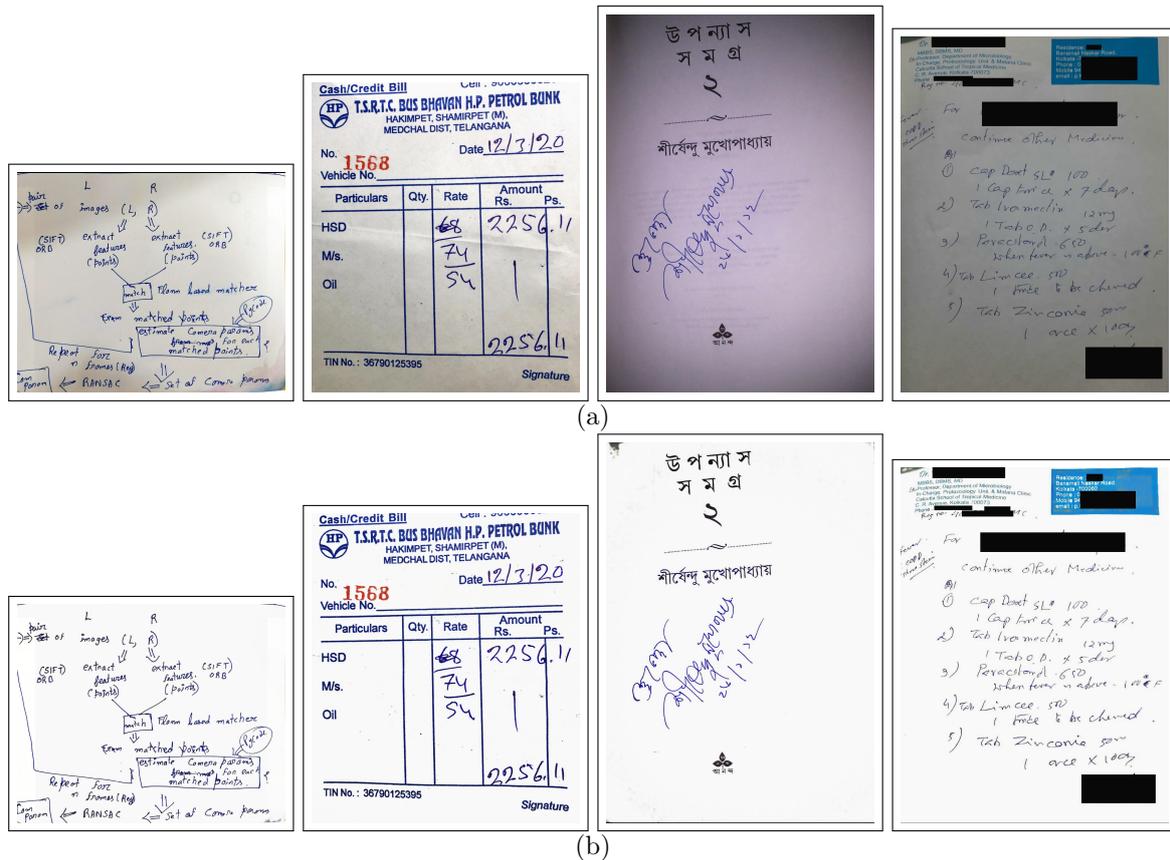
Figure 8: Typical examples of color clean up. (a) random inputs images; (b) cleaned outputs using M-32 based model on a mobile device

[8] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: 2016 IEEE Conf. on Computer Vision and Pattern Recognition (CVPR). pp. 770–778 (2016)

[9] He, S., Schomaker, L.: Deepotsu: Document enhancement and binarization using iterative deep learning. Pattern Recognition **91**, 379–390 (2019)

[10] Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H.: Mobilenets: Efficient convolutional neural networks for mobile vision applications (2017)

[11] Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: European Conf. on Computer Vision (ECCV) (2016)

[12] Kang, S., Iwana, B.K., Uchida, S.: Cascading modular u-nets for document image binarization. In: 2019 Int. Conf. on Document Analysis and Recognition (ICDAR). pp. 675–680 (2019)

[13] Kingma, P.D., Ba, L.J.: Adam: A method for stochastic optimization. Int. Conf. on learning representations (2015)

[14] Kise, K.: Page segmentation techniques in document analysis. In: Doermann, D., Tombre, K. (eds.) Handbook of Document Image Processing and Recognition, pp. 135–175. Springer London (2014)

[15] Kligler, N., Katz, S., Tal, A.: Document enhancement using visibility detection. In: 2018 IEEE/CVF Conf. on Computer Vision and Pattern Recognition. pp. 2374–2382 (2018)

[16] Krähenbühl, P., Koltun, V.: Efficient inference in fully connected crfs with gaussian edge potentials. In: Advances in Neural Information Processing Systems (2011)

[17] Krizhevsky, A.: Convolutional deep belief networks on cifar-10 (2010)

[18] Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems (2012)

[19] Lazzara, G., Géraud, T.: Efficient multiscale sauvola's binarization. Int. Journal of Document Analysis and Recognition **17**(2), 105–123 (2014)

[20] Li, X., Zhang, B., Liao, J., Sander, P.V.: Document rectification and illumination correction using a patch-based cnn. ACM Trans. on Graphics **38**(6) (Nov 2019)

[21] Lin, Y.H., Chen, W.C., Chuang, Y.Y.: Bedsr-net: A deep shadow removal network from a single document image. In: 2020 IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR). pp. 12902–12911 (2020)

[22] Liu, N., Zhang, D., Xu, X., Liu, W., Ke, D., Guo, L., Shi, S., Liu, H., Chen, L.: An iterative refinement framework for image document binarization with bhattacharyya similarity measure. In: 14th Int. Conf. on Document Analysis and Recognition. pp. 93–98. ICDAR '17, IEEE Computer Society (2017)

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

[24] Mondal, T., Coustaty, M., Gomez-Krämer, P., Ogier, J.: Learning free document image binarization based on fast fuzzy c-means clustering. In: 2019 Int. Conf. on Document Analysis and Recognition (ICDAR). pp. 1384–1389 (2019)

[25] Ntirogiannis, K., Gatos, B., Pratikakis, I.: Icfhr2014 competition on handwritten document image binarization (h-dibco 2014). In: 2014 14th Int. Conf. on Frontiers in Handwriting Recognition. pp. 809–813 (2014)

[26] Otsu, N.: A threshold selection method from gray-level histograms. IEEE Trans. on Systems, Man, and Cybernatics **9**(1), 62–66 (1979)

[27] Peng, X., Cao, H., Subramanian, K., Prasad, R., Natarajan, P.: Exploiting stroke orientation for crf based binarization of historical documents. In: 2013 12th Int. Conf. on Document Analysis and Recognition. pp. 1034–1038 (2013)

[28] Peng, X., Wang, C., Cao, H.: Document binarization via multi-resolutional attention model with drd loss. In: 2019 Int. Conf. on Document Analysis and Recognition (ICDAR). pp. 45–50 (2019)

[29] Pratikakis, I., Gatos, B., Ntirogiannis, K.: H-dibco 2010 - handwritten document image binarization competition. In: 2010 12th Int. Conf. on Frontiers in Handwriting Recognition. pp. 727–732 (2010)

[30] Pratikakis, I., Gatos, B., Ntirogiannis, K.: Icdar 2011 document image binarization contest (dibco 2011). In: 2011 Int. Conf. on Document Analysis and Recognition. pp. 1506–1510 (2011)

[31] Pratikakis, I., Gatos, B., Ntirogiannis, K.: Icfhr 2012 competition on handwritten document image binarization (h-dibco 2012). In: 2012 Int. Conf. on Frontiers in Handwriting Recognition. pp. 817–822 (2012)

[32] Pratikakis, I., Gatos, B., Ntirogiannis, K.: Icdar 2013 document image binarization contest (dibco 2013). In: 2013 12th Int. Conf. on Document Analysis and Recognition. pp. 1471–1476 (2013)

[33] Pratikakis, I., Zagori, K., Kaddas, P., Gatos, B.: Icfhr 2018 competition on handwritten document image binarization (h-dibco 2018). In: 2018 16th Int. Conf. on Frontiers in Handwriting Recognition (ICFHR). pp. 489–493 (2018)

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

[35] Pratikakis, I., Zagoris, K., Barlas, G., Gatos, B.: Icdar2017 competition on document image binarization (dibco 2017). In: 2017 14th IAPR Int. Conf. on Document Analysis and Recognition (ICDAR). vol. 01, pp. 1395–1403 (2017)

[36] Rad, M.S., Bozorgtabar, B., Marti, U., Basler, M., Ekenel, H.K., Thiran, J.: Srobb: Targeted perceptual loss for single image super-resolution. In: International Conference on Computer Vision (2019)

[37] Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention (2015)

[38] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: Imagenet large scale visual recognition challenge. Int. Journal of Computer Vision **115**(3), 211–252 (2015)

[39] S. Jung, Md. A. Hasan, C.K.: Water-filling: An efficient algorithm for digitized document shadow removal. In: 2018, 14th Asian Conf. on Computer Vision (ACCV) (2018)

[40] Sauvola, J., Pietikäinen, M.: Adaptive document image binarization. Pattern Recognition **33**, 225–236 (2000)

[41] Silva, J.M.M.D., Lins, R.D., Martins, F.M.J., Wachenchauzer, R.: A new and efficient algorithm to binarize document images removing back-to-front interference. Journal of Universal Computer Science **14**(2), 299–313 (2008)

[42] Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition (2014)

[43] Souibgui, M.A., Kessentini, Y.: De-gan: A conditional generative adversarial network for document enhancement. IEEE Trans. on Pattern Analysis and Machine Intelligence **early access**, 1–12 (2020)

[44] Tabatabaei, S.A., Bohlool, M.: A novel method for binarization of badly illuminated document images. 17th IEEE Int. Conf. on Image Processing pp. 3573–3576 (2010)

[45] Teichmann, M., Cipolla, R.: Convolutional crfs for semantic segmentation. In: British Machine Vision Conf. (2019)

[46] Tensmeyer, C., Martinez, T.: Document image binarization with fully convolutional neural networks. In: 14th Int. Conf. on Document Analysis and Recognition. pp. 99–104. ICDAR '17, IEEE Computer Society (2017)

[47] Valizadeh, M., Kabir, E.: An adaptive water flow model for binarization of degraded document images. Int. Journal on Document Analysis and Recognition pp. 1–12 (2012)

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

[49] Wang, B., Chen, C.L.P.: An effective background estimation method for shadows removal of document images. In: 2019 IEEE Int. Conf. on Image Processing (ICIP). pp. 3611–3615 (2019)

[50] Wang, J., Chuang, Y.: Shadow removal of text document images by estimating local and global background colors. In: ICASSP 2020 - 2020 IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP). pp. 1534–1538 (2020)

[51] Zamora-Martínez, F., España-Boquera, S., Castro-Bleda, M.J.: Behaviour-based clustering of neural networks applied to document enhancement. In: Sandoval, F., Prieto, A., Cabestany, J., Graña, M. (eds.) Computational and Ambient Intelligence. pp. 144–151. Springer Berlin Heidelberg, Berlin, Heidelberg (2007)

[52] Zhao, G., Liu, J., Jiang, J., Guan, H., Wen, J.: Skip-connected deep convolutional autoencoder for restoration of document images. In: 2018 24th Int. Conf. on Pattern Recognition (ICPR). pp. 2935–2940 (2018)

[53] Zhao, H., Gallo, O., Frosio, I., Kautz, J.: Loss functions for image restoration with neural networks. IEEE Transactions on Computational Imaging **3**(1), 47–57 (2017)