

Handwritten Text Segmentation via End-to-End Learning of Convolutional Neural Networks

1st Junho Jo

Electrical and Computer Engineering
INMC, Seoul National University
Seoul, Korea
jottue@ispl.snu.ac.kr

2nd Hyung Il Koo

Electrical and Computer Engineering
Ajou University
Suwon, Korea
hikoo@ajou.ac.kr

3rd Jae Woong Soh

Electrical and Computer Engineering
INMC, Seoul National University
Seoul, Korea
soh90815@ispl.snu.ac.kr

4th Nam Ik Cho

Electrical and Computer Engineering
INMC, Seoul National University
Seoul, Korea
nicho@snu.ac.kr

arXiv:1906.05229v1 [cs.CV] 12 Jun 2019

Abstract—We present a new handwritten text segmentation method by training a convolutional neural network (CNN) in an end-to-end manner. Many conventional methods addressed this problem by extracting connected components and then classifying them. However, this two-step approach has limitations when handwritten components and machine-printed parts are overlapping. Unlike conventional methods, we develop an end-to-end deep CNN for this problem, which does not need any preprocessing steps. Since there is no publicly available dataset for this goal and pixel-wise annotations are time-consuming and costly, we also propose a data synthesis algorithm that generates realistic training samples. For training our network, we develop a cross-entropy based loss function that addresses the imbalance problems. Experimental results on synthetic and real images show the effectiveness of the proposed method. Specifically, the proposed network has been trained solely on synthetic images, nevertheless the removal of handwritten text in real documents improves OCR performance from 71.13% to 92.50%, showing the generalization performance of our network and synthesized images.

Index Terms—handwritten text segmentation, text separation, data synthesis, class imbalance problem, optical character recognition

I. INTRODUCTION

Document digitization has been an important topic for the decades, and a huge number of methods have been proposed to address many kinds of sub-tasks such as optical character recognition (OCR), text-line segmentation, layout analysis, and so on [1]–[3]. Therefore, there have been many advances in machine-printed document understanding and handwritten text recognition. However, the understanding of mixed cases (*i.e.*, documents having handwritten and machine-printed texts on the same page) still remains a challenging problem, especially when they are overlapping. As shown in Fig. 1, these situations frequently occur in the formed documents, where we need to understand documents in the presence of handwritten notes and/or separate the handwritten and machine-printed texts.

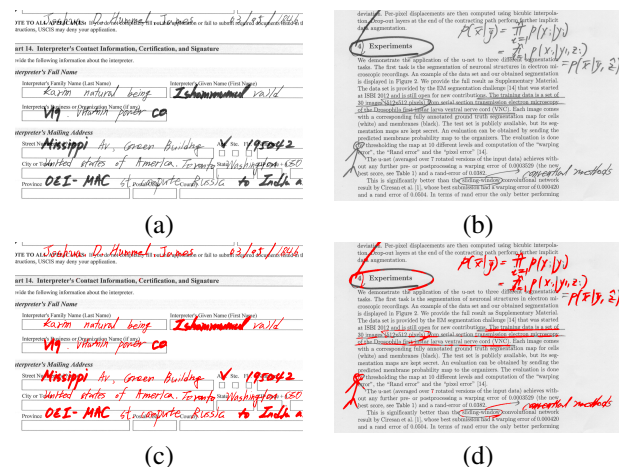


Fig. 1: The proposed method performs handwritten text segmentation from scanned document images: (a), (b) input images (synthetic and real), (c), (d) segmentation results (handwritten pixels are in red).

Many researchers addressed this problem by separating handwritten (or machine-printed) texts from the input [4]–[10]. In [8], they extracted connected components (CCs) and assigned feature vectors to them by exploiting the distribution of vectorized heights, widths, and distances between components. Finally, they classified each component by applying a k -nearest neighbor (NN) classifier. Similarly, Kandan *et al.* [6] classified each component by using support vector machines (SVMs) and k -NN classifiers. Also, they improved descriptors so that the algorithm is robust to deformations. Recently, CNNs outperform traditional methods based on hand-crafted features in a variety of applications [11]–[14], and Li *et al.* used CNNs to classify CCs [4]. Also, they incorporated conditional random fields into their framework to consider relations with neighboring CCs.

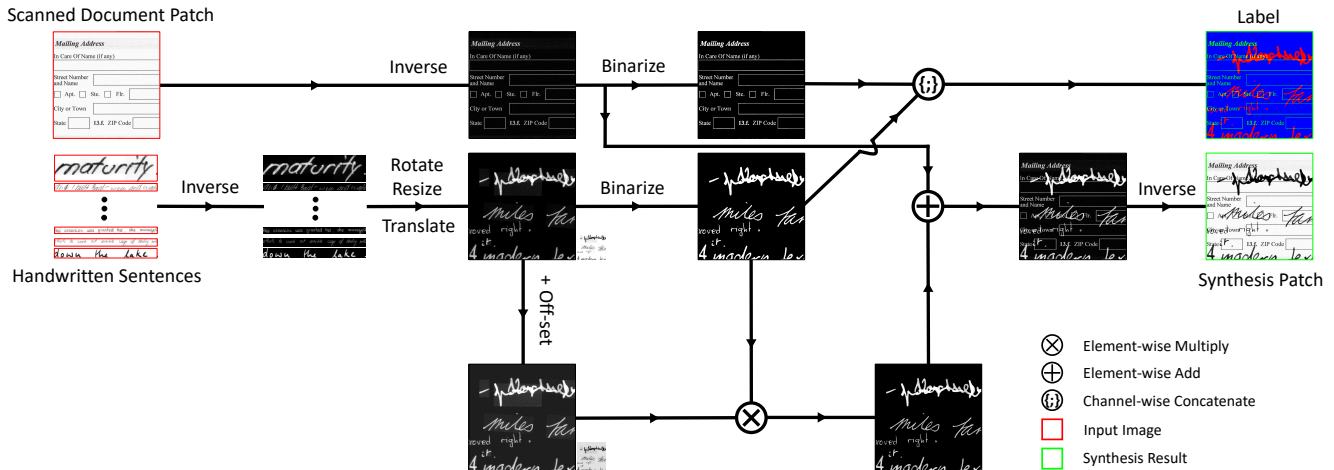


Fig. 2: Overall diagram of proposed data synthesis method.

However, most of these conventional methods employed binarization and CC extraction as essential preprocessing steps and thus used the binary classification of CCs. These two-step approaches have advantages in that they allow us to exploit a variety of conventional modules (e.g., binarization, CC extraction, etc.), which also means that they have drawbacks that the final performance heavily depends on the performance of each module. Also, the CC extraction methods are prone to errors when two different kinds of texts are overlapping.

To alleviate these problems, we propose a new handwritten text segmentation method based on an end-to-end CNN. To be precise, we formulate the task as a pixel-wise classification problem, assigning ‘+1’ for pixels of handwritten text and ‘-1’ for others (background, machine-printed text, table boundaries, and so on). For the segmentation network, we adopt the U-Net [13] that naturally exploits contextual information. In training the segmentation network, we address two challenges. First, the number of handwritten text pixels is much smaller than the number of other pixels (mainly due to background pixels) [15], so we develop a new loss function based on the conventional cross-entropy [11], [13]. Second, since there is no publicly available dataset (the manual pixel-level annotation of documents is time-consuming and costly), we also propose a new synthesis method that can generate realistic images along with ground truth labels.

The contributions of this paper can be summarized as follows:

- To the best of our knowledge, this is the first method that applies end-to-end learning to the separation of mixed-handwritten-machine-printed texts.
- For training the network under imbalanced situations, we propose a new loss function based on cross-entropy.
- For training, we also develop a data synthesis method, yielding realistic scanned documents as shown in Fig. 3(b).

II. DATASET SYNTHESIS

Although deep learning methods outperform conventional methods in many fields, training the deep networks requires a huge number of training samples. Especially, for the learning of segmentation networks (our application), pixel-level annotations are needed. However, there is no publicly available dataset for this goal, and we address this problem by synthesizing training samples

A. Scan image dataset

Synthesizing realistic images from scratch is a difficult task, and we develop a method that uses existing scanned images. We use 13,353 sentence images of IAM dataset [16] as the handwritten texts, which was written by a variety of writers. For machine-printed parts, PRImA dataset [17] is used, which consists of scanned images of magazines and technical/scientific publications. Additionally, we manually crawled 141 images of questionnaire forms from the Internet. We augmented dataset by including these images in training dataset so that our dataset covers a wide range of formed documents. Typical examples are shown in *red* boxes in Fig. 2.

B. Dataset Synthesis

For the realistic data synthesis, our main considerations are preserving textures of handwritten text images and noise of original documents. First, it is noted that textures of handwritten text images can be crucial evidence to differentiate them from the machine printed texts. Secondly, consistent noise inherited from the scanning process must be preserved to diminish discrepancies between the distributions of synthetic and real data. To be precise, if we simply add a handwritten text image and a machine-printed text image, then the background will be saturated, and most of the scan-noises will disappear. We address this issue by inverting their intensities because backgrounds do not suffer from saturation if two inverted images are added. Another issue is undesirable block artifacts shown in Fig. 3(a). Actually, these artifacts are from IAM

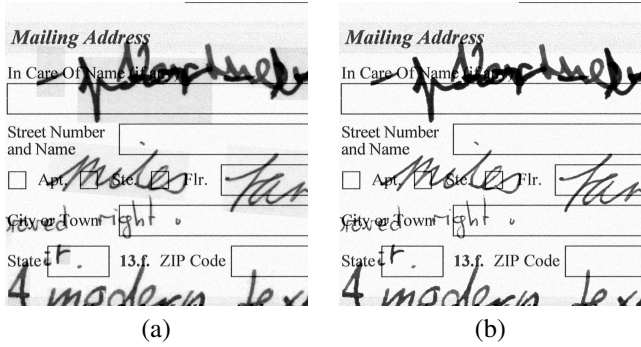


Fig. 3: (a) synthesis result w/o artifacts handling, (b) synthesis result using the proposed method.

dataset, since sentence images in IAM dataset were made by concatenating separate word images. To remove these artifacts, we extract only handwritten text pixels by multiplying images with its binary mask generated by *Otsu* binarization method [19]. As shown in Fig. 3(b), the proposed method yields realistic images, which have spatially consistent scan-noise.

In order to reflect the diversity of appearance observed in real environments, we apply randomized transformations, such as resizing, translation and rotation to each handwritten sentence image. We also augment the dataset by adding random off-set values to the handwritten text patch to simulate a variety of intensities of handwritten texts. We have synthesized 146,391 patches for training and 8,128 for validation. The overall synthesis method is shown in Fig. 2. We will make our synthesized dataset publicly available.

III. PROPOSED METHOD

With synthesized training samples, we train a network in an end-to-end manner. This section presents our network structure and loss functions that are appropriately designed for our purpose and environment.

A. Network Structure

As a neural network architecture, we adopt U-Net [13] that consists of encoder and decoder parts as shown in Fig. 4. The encoder captures the context of images by using a large-scale receptive field (downsampling operators), and the decoder generates high-resolution segmentation results by using contextual information and features from the encoder. As downsampling operators, we empirically select max-pooling instead of strided-convolutions, and 4×4 transposed-convolution layers with stride 2 are used to up-sample the concatenation of encoder and decoder signals.

B. Cross entropy based loss function

The cross-entropy loss function is commonly used for training the recognition [12], [14] and segmentation networks [11], [13], which is described as

$$\mathcal{L}_{\text{CE}}(\theta) = \frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C \text{CE}(n, c), \quad (1)$$

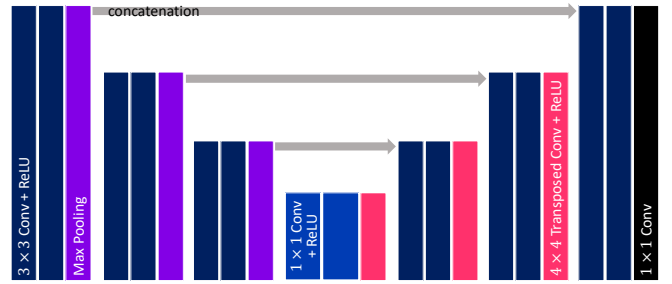


Fig. 4: Our network architecture based on U-Net.

where

$$\text{CE}(n, c) = -t_{n,c} \log y_{n,c}, \quad (2)$$

where $n \in \{1, 2, \dots, N\}$ is the pixel index and $c \in \{1, 2, \dots, C\}$ is the class index. In our case, C is set to 2, since our model makes a decision just whether this pixel is handwritten text pixel or not. Also, $t_{n,c} \in \{0, 1\}$ and $y_{n,c} \in [0, 1]$ represent one-hot encoding of the ground truth label and the softmax result of the network, respectively. The θ denotes parameters of the network.

However, eq.(1) is likely to yield a poor local minimum for our task. That is, in most document images, the number of background pixels is approximately 20 times larger than that of text pixels. Therefore, the model is likely to converge to a trivial solution that classifies all pixels as background (the *class imbalance* problem). Moreover, background pixels consist of many easy cases and a tiny number of hard cases. In other words, most background pixels can be easily classified even by a simple thresholding method, and the CNN is very likely to converge to a sub-optimal solution by focusing on easy but dominant cases. That is, the loss summed over a large number of easy background examples would *overwhelm* the loss of rare hard examples, i.e., “*many a little makes a mickle*” (the *overwhelming* problem).

C. Dynamically Balanced Cross Entropy

In order to alleviate the aforementioned *class imbalance* problem, we propose a new loss function:

$$\mathcal{L}_{\text{DBCE}}(\theta) = \frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C \frac{1}{\beta(c) + \epsilon} \text{CE}(n, c), \quad (3)$$

where

$$\beta(c) = \frac{1}{N} \sum_{n=1}^N [t_{n,c} == 1]. \quad (4)$$

Unlike eq.(1), $\text{CE}(n, c)$ is (dynamically) divided by the *frequency* of pixels in each class (in *mini-batch*) in $\mathcal{L}_{\text{DBCE}}(\theta)$. That is, the amount of contribution of each pixel is weighted by the scarcity of its class (fewer cases will have larger weights). By employing our loss function, the model can be trained even in imbalanced situations.

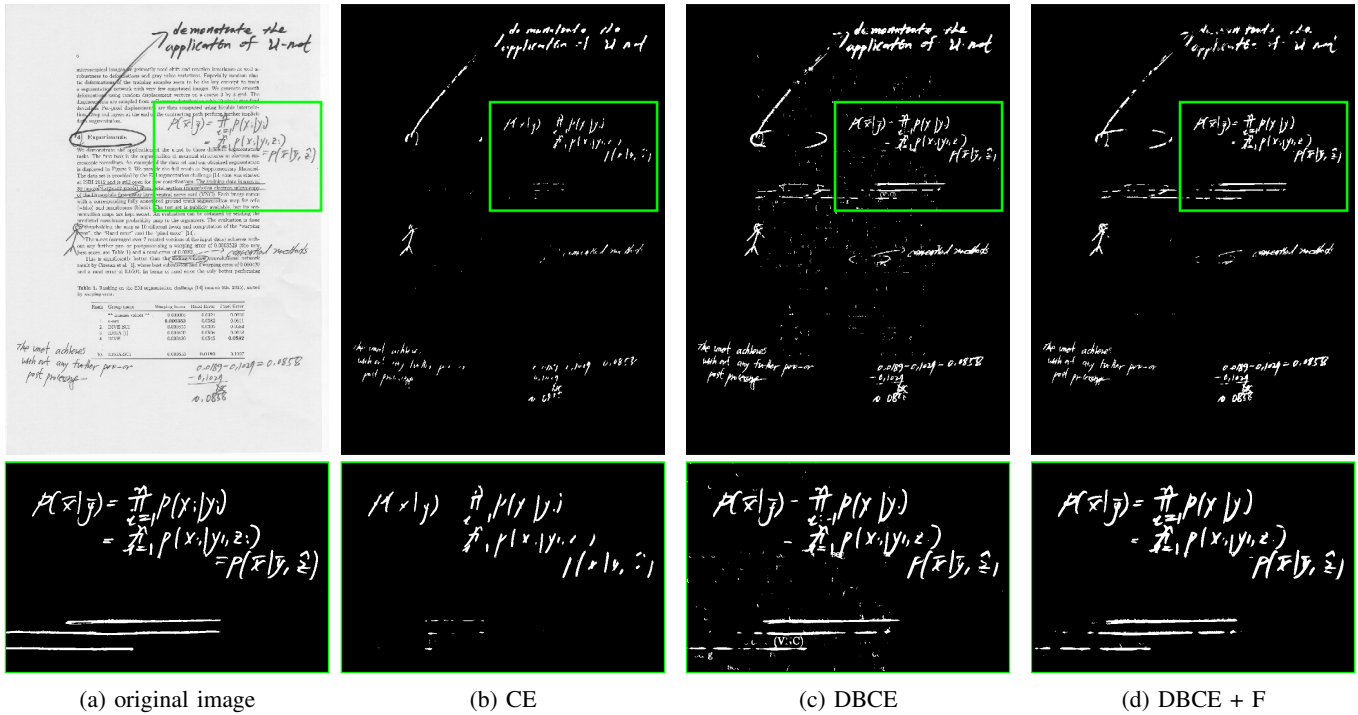


Fig. 5: Segmentation results on a real scribbled document by the proposed network trained by each loss function. CE: conventional cross entropy, DBCE: dynamically balanced cross entropy, F: focal loss.

D. Focal Loss

In order to alleviate the *overwhelming* problem, we adopt focal loss in [14]. For the presentation of focal loss, we first define $\mathbf{FCE}(n, c)$ as

$$\mathbf{FCE}(n, c) = - (1 - y_{n,c})^\gamma t_{n,c} \log y_{n,c}, \quad (5)$$

where γ is the hyperparameter that determines the boost degree of the penalty. As shown in eq.(5), the term $\mathbf{FCE}(n, c)$ is the $(1 - y_{n,c})^\gamma$ scaled version of $\mathbf{CE}(n, c)$ in eq.(2). This scaling factor automatically lessens the contribution of easy examples and makes the model focus on hard examples during the training. By putting two ideas together, the final loss function is given by

$$\mathcal{L}_{\text{DBCEF}}(\theta) = \frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C \frac{1}{\beta(c) + \epsilon} \mathbf{FCE}(n, c). \quad (6)$$

E. Training details

We have trained the network using Adam optimizer [18] with a mini-batch size of 32. We used 0.0002 as the initial learning rate with 0.8 decay rate in every 30 epochs. For hyperparameter of loss function, we empirically set $\epsilon = 0.0001$ and $\gamma = 1$.

IV. EXPERIMENTAL RESULTS

In this section, we perform an ablation study to see whether our loss functions are effective for training the network. Then, we will show the performance of the proposed method on synthetic and real data. We did not perform the comparison

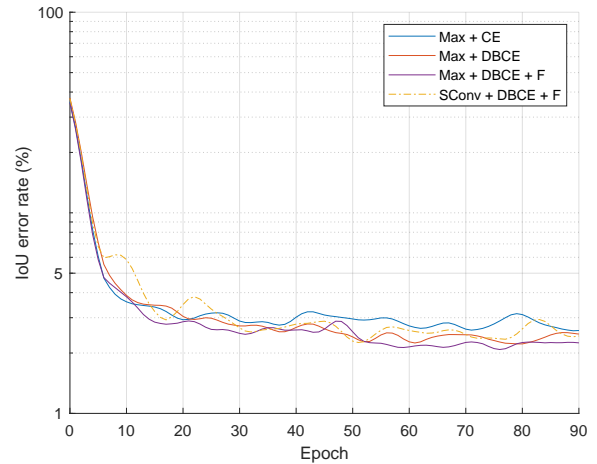


Fig. 6: IoU error rate of handwritten text in validation set. *Max*: max-pooling, *SConv*: strided-convolution, *CE*: conventional cross entropy, *DBCE*: dynamically balanced cross entropy, *F*: focal loss.

with existing works [4]–[10], since there is none that publicly provides the code and data to compare the performance. In a recent research of [4], they tested CC-level and region-level segmentation with their own *TestPaper 1.0* dataset, and the *Maurdor* dataset from [20]. However, these datasets are currently not accessible, and also they cannot be directly compared with ours because we deal with pixel-level results.

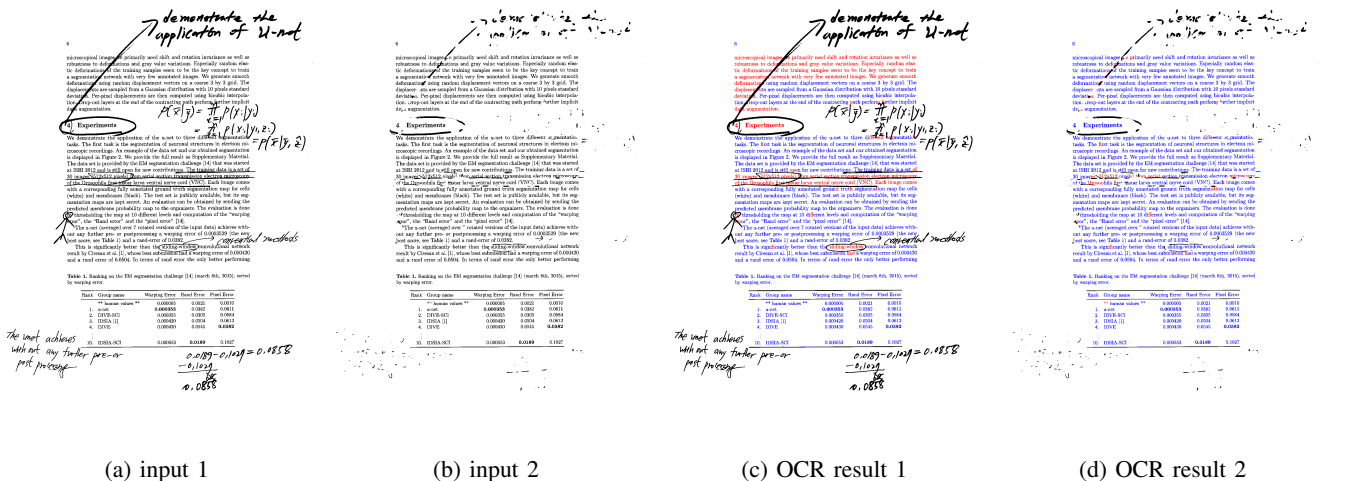


Fig. 7: Comparisons of OCR performance. (a) Input 1 to OCR (real scribbled image), (b) Input 2 to OCR (removal of handwritten pixels from (a)), (c) and (d) OCR results for input 1 and 2 (blue: correctly recognized characters, red: missing or incorrect ones).

Hence, we will instead make our dataset and codes publicly available for future research and comparisons.

A. Ablation Study

Fig. 5 demonstrates the effect of each loss function on the segmentation results, where the first column shows the input and the ground truth of segmentation results on green-box region. As shown in Fig. 5(b), using the conventional cross entropy loss function (\mathcal{L}_{CE}) for training, a lot of handwritten pixels are classified as background, which means that the trained model experiences the *class imbalance* problem. By applying the proposed \mathcal{L}_{DBCE} , we can see that the *class imbalance* problem is quite mitigated as shown in Fig. 5(c). However, there are lots of misclassified pixels in the machine printed text region, which is the example of *overwhelming* problem. To alleviate this problem, we incorporate focal loss [14] with proposed \mathcal{L}_{DBCE} . Finally, we can get the segmentation results without any degradation from *imbalance* situations as shown in Fig. 5(d).

For the quantitative comparisons, we evaluate the proposed method using the pixel-level intersection-over-union (IoU) on synthesized sets. Fig. 6 shows that the proposed loss function (\mathcal{L}_{DBCEF}) improves validation performance. Also, as shown in Table I, \mathcal{L}_{DBCE} and focal loss term achieves 1.32%p and 0.60%p improvements of test performance on synthesized set, respectively. These qualitative and quantitative experimental results showed that \mathcal{L}_{DBCE} and focal loss term is meaningfully functioning during network training, *i.e.*, mitigating *class imbalance* and *overwhelming* problem well.

B. Generalization performance

In the case of *real* test-set, IoU evaluation is infeasible due to the lack of (pixel-level) ground truth. Rather, we measure the OCR performance on handwritten-pixel-removed-images, which is naturally proportional to the handwritten text segmentation performance. To be precise, given a scribbled image

TABLE I: IoU results on synthesized test set. The best results are highlighted in **bold face** and the second best results are underlined. *H*: handwritten text, *Max*: max-pooling, *SConv*: strided-convolution, *CE*: conventional cross entropy, *DBCE*: dynamically balanced cross entropy, *F*: focal loss

Model	Number of Parameters (M)	IoU (%)	
		non-H	H
<i>Max</i> + <i>CE</i>	6.61	99.89	95.88
<i>Max</i> + <i>DBCE</i>	6.61	<u>99.93</u>	<u>97.20</u>
<i>Max</i> + <i>DBCE</i> + <i>F</i>	6.61	99.94	97.80
<i>SConv</i> + <i>DBCE</i> + <i>F</i>	7.39	99.92	97.11

TABLE II: OCR performance on a real scribbled-document image. Accuracy is calculated by Correct / (Correct + Incorrect + Missing). Visualized results are shown in Fig. 7.

	Document states		
	original	scribbled	separated
Correct	2,141	1,584	2,071
Incorrect	3	125	164
Missing	7	518	4
Accuracy (%)	99.54	71.13	92.50

like Fig. 7(a), we evaluate the OCR performances of original documents (w/o handwritten components) and handwritten-pixel-removed-images (Fig. 7(b)). As shown in Table II and Fig. 7, there are a lot of missing or incorrectly detected characters in the scribbled document, mainly due to scribbles overlapping the machine-printed text. After removing them, which is segmented as handwritten text by proposed network, the OCR performance is improved from 71.13% to 92.50%. Note that the model is trained only with the synthesized data,

and these results show that the model has learned features having generalization ability.

V. CONCLUSION

In this paper, we have proposed a method to separate handwritten text from the machine-printed documents based on a deep neural network. Unlike conventional methods, we proposed a method that works in an end-to-end manner, and addressed the *class imbalance* and *overwhelming* problems in the training phase. The network was trained with synthetic training samples generated by the proposed synthesis method. The experimental results show that the proposed model also works well for real document images. Although the proposed method shows a good handwritten text extraction performance, it can reconstruct only handwritten part when it is overlapping with the machine-printed text. As a future work, we will address the layer separation problem to reconstruct both components.

REFERENCES

- [1] Ray Smith, An overview of the tesseract ocr engine, in Document Analysis and Recognition, 2007. ICDAR 2007. Ninth International Conference on. IEEE, 2007, vol. 2, pp. 629633.
- [2] Jewoong Ryu, Hyung Il Koo, and Nam Ik Cho, Language-independent text-line extraction algorithm for handwritten documents, IEEE Signal processing letters, vol. 21, no. 9, pp. 11151119, 2014.
- [3] Hyung Il Koo and Nam Ik Cho, Text-line extraction in handwritten Chinese documents based on an energy minimization framework, IEEE Transactions on Image Processing, vol. 21, no. 3, pp. 11691175, 2012.
- [4] Xiao-Hui Li, Fei Yin, and Cheng-Lin Liu, Printed/handwritten texts and graphics separation in complex documents using conditional random fields, in 2018 13th IAPR International Workshop on Document Analysis Systems (DAS). IEEE, 2018, pp. 145150.
- [5] Mathias Seuret, Marcus Liwicki, and Rolf Ingold, Pixel level handwritten and printed content discrimination in scanned documents, in Frontiers in Handwriting Recognition (ICFHR), 2014 14th International Conference on. IEEE, 2014, pp. 423428.
- [6] R Kandan, Nirup Kumar Reddy, KR Arvind, and AG Ramakrishnan, A robust two level classification algorithm for text localization in documents, in International Symposium on Visual Computing. Springer, 2007, pp. 96105.
- [7] Xujun Peng, Srirangaraj Setlur, Venu Govindaraju, and Ramachandrala Sitaram, Handwritten text separation from annotated machine printed documents using markov random fields, International Journal on Document Analysis and Recognition (IJ DAR), vol. 16, no. 1, pp. 116, 2013.
- [8] Jurgen Franke and Matthias Oberlander, Writing style detection by statistical combination of classifiers in form reader applications, in Document Analysis and Recognition, 1993., Proceedings of the Second International Conference on. IEEE, 1993, pp. 581584.
- [9] Konstantinos Zagoris, Ioannis Pratikakis, Apostolos Antonacopoulos, Basilis Gatos, and Nikos Papamarkos, Handwritten and machine printed text separation in document images using the bag of visual words paradigm, in Frontiers in Handwriting Recognition (ICFHR), 2012 International Conference on. IEEE, 2012, pp. 103108.
- [10] Abdel Belad, KC Santosh, and Vincent Poulain dAndecy, Handwritten and printed text separation in real document, arXiv preprint arXiv:1303.4614, 2013.
- [11] Jonathan Long, Evan Shelhamer, and Trevor Darrell, Fully convolutional networks for semantic segmentation, in Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 3431 3440.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, Deep residual learning for image recognition, in Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770778.
- [13] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, U-net: Convolutional networks for biomedical image segmentation, in International Conference on Medical image computing and computer-assisted intervention. Springer, 2015, pp. 234241.
- [14] Tsung-Yi Lin, Priyal Goyal, Ross Girshick, Kaiming He, and Piotr Dollar, Focal loss for dense object detection, IEEE transactions on pattern analysis and machine intelligence, 2018.
- [15] Zhi-Hua Zhou and Xu-Ying Liu, Training costsensitive neural networks with methods addressing the class imbalance problem, IEEE Transactions on Knowledge and Data Engineering, vol. 18, no. 1, pp. 6377, 2006.
- [16] U-V Marti and Horst Bunke, The iam-database: an english sentence database for offline handwriting recognition, International Journal on Document Analysis and Recognition, vol. 5, no. 1, pp. 3946, 2002.
- [17] Apostolos Antonacopoulos, David Bridson, Christos Papadopoulos, and Stefan Pletschacher, A realistic dataset for performance evaluation of document layout analysis, in Document Analysis and Recognition, 2009. ICDAR09. 10th International Conference on. IEEE, 2009, pp. 296300.
- [18] Diederik P Kingma and Jimmy Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980, 2014.
- [19] Otsu, Nobuyuki. "A threshold selection method from gray-level histograms." IEEE transactions on systems, man, and cybernetics 9.1 (1979): 62-66.
- [20] Brunessaux, Sylvie, et al. "The maudor project: Improving automatic processing of digital documents." 2014 11th IAPR International Workshop on Document Analysis Systems. IEEE, 2014.