

LS-HDIB: A LARGE SCALE HANDWRITTEN DOCUMENT IMAGE BINARIZATION DATASET

Kaustubh Sadekar Ashish Tiwari Prajwal Singh Shanmuganathan Raman

CVIG Lab, Indian Institute of Technology Gandhinagar
{ sadelkar.k, ashish.tiwari, singh_prajwal, shanmuga }@iitgn.ac.in

ABSTRACT

Handwritten document image binarization is challenging due to high variability in the written content and complex background attributes such as page style, paper quality, stains, shadow gradients, and non-uniform illumination. While the traditional thresholding methods do not effectively generalize on such challenging real-world scenarios, deep learning-based methods have performed relatively well when provided with sufficient training data. However, the existing datasets are limited in size and diversity. This work proposes LS-HDIB - a large-scale handwritten document image binarization dataset containing over a million document images that span numerous real-world scenarios. Additionally, we introduce a novel technique that uses a combination of adaptive thresholding and seamless cloning methods to create the dataset with accurate ground truths. Through an extensive quantitative and qualitative evaluation over eight different deep learning based models, we demonstrate the enhancement in the performance of these models when trained on the LS-HDIB dataset and tested on unseen images.

Index Terms— Document Image Binarization, Deep Learning, Adaptive Thresholding.

1. INTRODUCTION

Handwritten document image binarization is generally modeled as a classification problem in which intra-image pixels are assigned to either of the two classes: handwritten content (the foreground) or the background. Document image binarization has been an active research area for decades owing to its importance as an essential pre-processing step in facilitating several document image processing tasks such as optical character recognition, handwriting matching, document translation, document summarization, and changing the background. Handwritten documents range from ancient documents, old legal records, and ledgers to music scores and handwritten bills. These documents often degrade over time and become difficult to comprehend. Common degradation scenarios include crumpled pages, poor foreground-background contrast, stains, paper aging, faded characters, uneven illumination, and bleed/show-through. Furthermore, handwritten documents possess various page styles, i.e., grids, lines, staff annotation styles, and partially blank pages that increase the difficulty of segmenting the foreground content. The variability in the type and the thickness of strokes also increases the complexity. These challenges make handwritten document image binarization extremely difficult. Traditional methods like Otsu image segmentation [1] and adaptive thresholding [2] fail to address the aforementioned challenges completely. Since these algorithms utilize only the low-level features, they fail to capture the

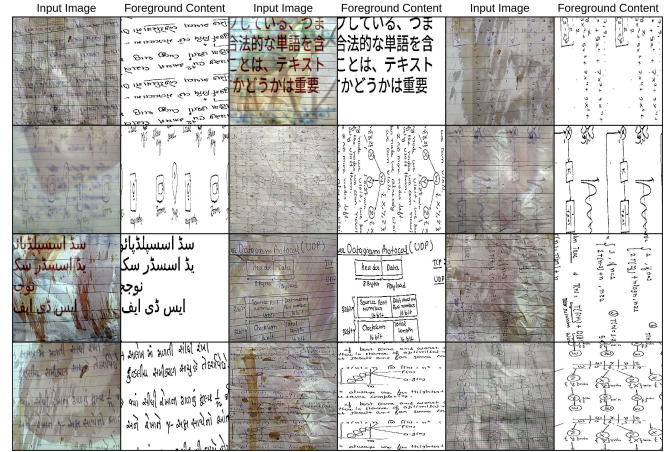


Fig. 1. A few samples of handwritten document images obtained from the proposed LS-HDIB dataset.

wide range of variability inherent in the handwritten documents, limiting their ability to distinguish the background from the foreground content. The high-level features can differentiate text pixels from background noises handling the degradations better. However, using them solely can cause loss of low-level information like character edges and contours, making it insufficient to address the binarization problem. Thus far, deep learning methods have shown promising results in segmenting foreground and background content [3, 4] by incorporating both low-level and high-level image features. Deep learning models rely heavily on a large amount of data for better generalization [5]. However, document image binarization lacks such a large and diverse dataset that covers numerous real-world scenarios. While the focus has been on designing new robust networks, little attention has been given to scale up the existing datasets.

To address the aforementioned challenges and to generate a large scalable dataset, we propose a Large Scale Handwritten Document Image Binarization dataset (LS-HDIB) that contains over a million handwritten document images. We propose a simple and effective method for generating the LS-HDIB dataset with accurate ground truths containing segmented handwritten content (see Fig. 1). Interestingly, the proposed method requires no manual intervention for generating these segmented ground truths. The primary **contributions** of this work are: (i) A large scale dataset (LS-HDIB) containing over a million images with accurate ground truths for handwritten document image binarization. (ii) A scalable and efficient method to generate and extend the proposed dataset.

¹This work is supported by SERB IMPRINT-2 grant.

2. RELATED WORK

The standard approaches for document image binarization are classified into (i) *global methods* [1, 6] which use a single threshold value and (ii) *local methods* [7] which use adaptive threshold values for separating the foreground and the background content. While the global methods fail to handle complex degradations, the local methods are computationally expensive and driven mainly by manual parameter tuning. Several deep neural network architectures have been designed for document image binarization [8, 9, 10]. Researchers have used fully convolutional neural networks [11, 12], recurrent neural networks [13], encoder-decoder frameworks [9, 12], and generative adversarial networks [11, 14] to address document binarization. While the performance of these learning-based frameworks depends on the robustness of the network design, another aspect critical to their performance is the size and the versatility of the training data. The existing datasets [15, 16, 17, 18, 19, 20, 21, 22, 23, 24] do not completely span complex degradations, page styles, and illumination variations. Owing to a very large space spanned by the possible handwritten content and background variations, there is a need to develop a method to create a large and diverse dataset for handwritten document image binarization. In this work, we propose a simple yet effective method to create a large-sized dataset that can potentially circumvent the aforementioned limitations.

3. METHOD

3.1. Dataset Generation

We propose a novel data generation technique that uses a combination of adaptive thresholding [2] and mixed gradient seamless cloning [25], as described in Fig. 2. We collect the images of a variety of handwritten content over a plain background and refer to these images as *full-length document images* (I_{doc}). We collected over 450 full-length document images. Around 400 images contain handwritten content in various forms such as alpha-numeric characters, electrical circuit diagrams, control system schematics, chemical molecular structures, and flow charts that are not present in the existing datasets. Further, to diversify the types, thickness, and styles of strokes, we obtained these images from 21 different persons. Around 50 document images were digitally created, using Google Translate, in 13 different languages, including English, Urdu, Mandarin, Portuguese, Russian, French, Hindi, Telugu, Malayalam, Punjabi, Gujarati, Japanese, and Korean in 21 different font styles of various sizes and colors.

We apply the *adaptive thresholding* (\mathcal{T}) [2] on (I_{doc}) to obtain the *segmented ground truth images* (I_{gt}) such that $I_{gt} = \mathcal{T}(I_{ref})$. Next, we generate a total of 10,944 unique *content images* (I_c) by cropping and augmenting multiple patches of size 480×480 from the full-length document images (I_{doc}). Figure 3(a) shows a few sample images containing a variety of written content in different languages and font styles, including diagrams and texts from different subject domains. While the written content (foreground) associated with the generated ground truth remains unchanged, the background can vary depending on different page styles and degradation scenarios. We obtain multiple document images by merely changing the background with essentially the same ground truth. Moreover, this method can automatically generate ground truths saving hours of tedious manual annotation. To generate different backgrounds, we manually capture pages with a wide variety of *page styles* (I_p) and *degradation effects* (I_d) that are predominant in the real-world. We then use mixed gradient-based seamless cloning (\mathcal{C}) [25] to blend

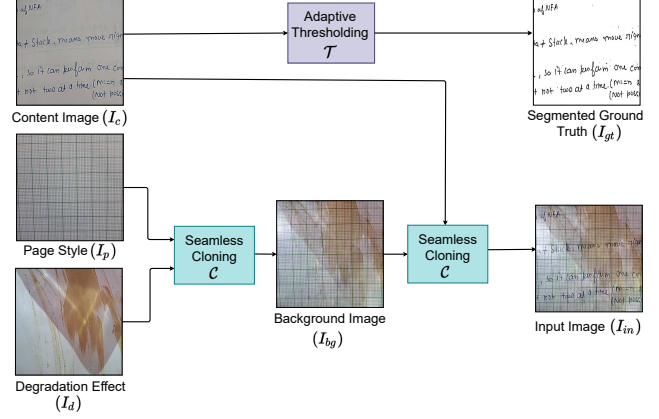


Fig. 2. Block schematic of the proposed method for generating LS-HDIB dataset.

multiple patches of (I_p) and (I_d) to generate 20,484 photorealistic *background images* (I_{bg}) such that $I_{bg} = \mathcal{C}(I_p, I_d)$.

For a better visual understanding, we show multiple example images of different page styles and degradation effects in Fig. 3(a) and Fig. 3(b), respectively. Finally, we combine the content images (I_c) and the background images (I_{bg}) again using seamless cloning (\mathcal{C}) to generate the handwritten document images (I_{in}) such that $I_{in} = \mathcal{C}(I_c, I_{bg})$. To generate the LS-HDIB dataset, we randomly sample 100 background images (I_{bg}) for each of the 10,994 handwritten content images (I_c). In this way, we obtain a total of 1.09 million images in the LS-HDIB dataset. It is important to note that we can scale up the dataset size for up to 200 times by considering all the background images instead of just 100. With adaptive thresholding and mixed gradient-based seamless cloning, we have been able to generate accurate ground truths without the requirement of any manual annotation and generate a wide variety of degraded handwritten document images. In addition to the inherent scalability and ease of ground truth generation, the proposed dataset generation method is computationally less expensive compared to the other deep-learning-based generative methods such as the one proposed in [14]. As evident from Table 1, most of the publicly available datasets contain the written content only over plain pages. In contrast, the LS-HDIB dataset includes images with various page styles like ruled lines, gridlines, and partially blank pages that are evident in our day-to-day encounters. Further, they lack the document images with realistic degradations like crumpled pages, non-uniform illumination, and shadow gradients that are well incorporated in the proposed LS-HDIB dataset. These attributes enhance the diversity and the versatility of the proposed dataset.

3.2. Training Details

We use eight widely used deep networks - DeepLabV3 [26], DeepLabV3+ [27], Feature Pyramid Networks (FPN) [28], LinkNet [29], Multi-scale Attention net (MANet) [30], Pyramid Attention Network (PAN) [31], Pyramid Scene Parsing Network (PSPN) [32], and U-Net [33] - to understand the effectiveness of the proposed dataset for the handwritten document image binarization task. We have carefully chosen these networks as they collectively span the various deep learning based approaches [11, 12, 9, 14] that have been adopted thus far for the binarization task.

Each of the eight networks is trained under three different training regimes. We follow the standard train, validation, and test split

Datasets	Page Styles						Degradation Effects								
	Uniform ruled lines	Non-uniform ruled lines	Grid lines	Staff notation lines	Partially blank pages	Plain page	Shadow gradients	Oily patches	Ink bleed-through	crumpled pages	Non-uniform illumination	Noisy background	Liquid stains	Poor foreground-background contrast	Punched, stapled or torn pages
DIBCO09	x	x	x	x	✓	✓	x	x	✓	x	x	x	✓	✓	x
HDIBCO10	✓	x	x	x	x	✓	x	✓	✓	x	x	✓	✓	✓	x
DIBCO11	x	x	x	x	x	✓	x	x	✓	x	x	✓	✓	✓	x
HDIBCO12	✓	x	x	x	x	✓	x	✓	✓	x	✓	✓	✓	✓	x
DIBCO13	x	x	x	x	x	✓	x	✓	✓	✓	x	✓	✓	✓	x
HDIBCO14	x	x	x	x	x	✓	x	x	✓	x	x	x	✓	✓	x
PHIBD12	✓	x	x	x	x	✓	x	✓	✓	x	x	✓	✓	✓	✓
HDIBCO16	x	x	x	x	x	✓	x	✓	✓	x	✓	✓	✓	✓	x
DIBCO17	x	x	x	x	x	✓	x	✓	✓	x	x	✓	✓	✓	x
DIBCO18	x	x	x	x	x	✓	x	✓	✓	x	x	✓	✓	✓	x
Bickley Diary	x	x	x	x	✓	✓	✓	✓	x	x	x	✓	✓	✓	✓
Palm Leaf Manuscript	x	x	x	x	✓	✓	✓	✓	x	x	x	✓	✓	✓	✓
LS-HDIB (Ours)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	x

Table 1. Comparison of page styles and degradation effects available across different publicly available datasets and LS-HDIB dataset.

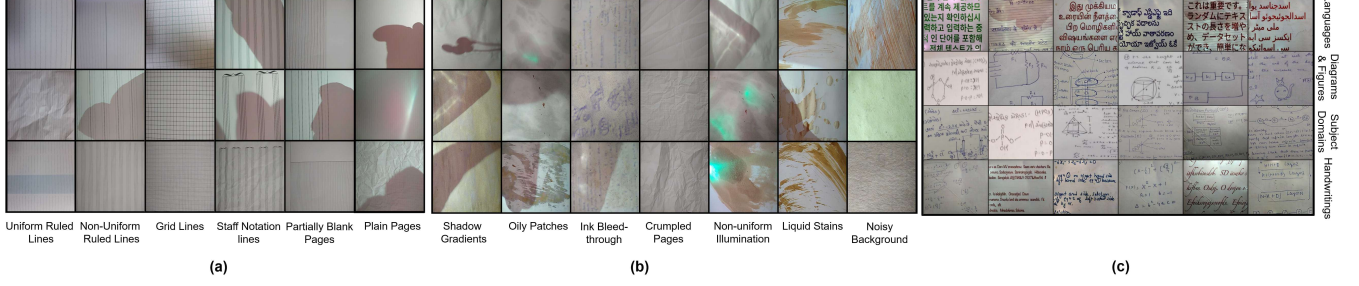


Fig. 3. A few sample images depicting different page styles available in LS-HDIB dataset.

of 80%, 10%, and 10%, respectively, for each regime.

(i) **Regime 1:** The deep models are trained only on the *baseline dataset* obtained by combining ten different publicly available datasets: DIBCO09 [15], HDIBCO10 [16], DIBCO11 [17], HDIBCO12 [18], DIBCO13 [19], HDIBCO14 [20], PHIBD12 [24], HDIBCO16 [21], DIBCO17 [22], DIBCO18 [23]. However, even after combining ten different datasets, the size of the baseline dataset is relatively small (order of 1000 images). Therefore, we crop the full-length document images of the baseline dataset to the size 480×480 with a stride of 240 pixels and perform rotation (90° , 180° , and 270°) and horizontal flip to obtain around 6000 images, 5000 for training and 1000 for testing. While the size of the baseline dataset can further be increased by reducing the stride value, this leads to greater overlap and high redundancy in the foreground content across different content images.

(ii) **Regime 2:** Each deep model is trained only on the LS-HDIB training set. Although the LS-HDIB dataset has over 1 million images, we use only 5000 images for training to have a fair performance comparison.

(iii) **Regime 3:** We combine both the LS-HDIB and the baseline dataset, by randomly selecting 2500 images from each dataset to train all the deep models on a total of 5000 images. Regime 3 is targeted towards establishing the efficacy of augmenting the proposed dataset to the existing ones.

The models are trained for a maximum of 20 epochs with learning rate of 0.01 using the Adam optimizer with default parameters. The training is performed on the NVIDIA RTX 2080 Ti GPU with the batch size of 8.

Loss Function. We use binary cross-entropy loss to train the segmentation models, as described in Equation 1.

$$\mathcal{L} = -I_{gt} \log(\mathcal{F}(I_{in})) - (1 - I_{gt}) \log(1 - \mathcal{F}(I_{in})) \quad (1)$$

Here, \mathcal{F} represents the functional form of deep model. Binary cross-entropy loss is found to be more effective than MSE loss for classification tasks [34].

4. EXPERIMENTAL ANALYSIS

We demonstrate the effectiveness of the proposed dataset for handwritten document image binarization through an extensive quantitative and qualitative analysis. We compare the performance of different deep models trained to observe how well the LS-HDIB dataset enhances the generalization capability of the deep models.

We use four different datasets containing challenging scenarios for testing the deep models: the Bickley Diary dataset [35], the Palm Leaf Manuscript dataset [36], the DIBCO test set, and the LS-HDIB test set. We use three popular metrics to evaluate the network performance trained under different regimes: F-measure (F_{score}) [37], pseudo F-measure (PF_{score}) [37], and Peak Signal to Noise Ratio (PSNR) [37].

As shown in Table 2, the F_{score} , PF_{score} , and PSNR is maximum over all the deep models trained under Regime 2 for LS-HDIB and the Bickley Diary dataset. Further, the performance under Regime 3 is better than that of Regime 1, indicating that augmenting the LS-HDIB dataset to the baseline dataset enhances network performance. Qualitatively, the foreground content of the image affected by liquid stains, noisy background, and poor foreground to background contrast is well recovered under Regime 2, as shown in Fig. 4(a) and 4(b). For the Palm Leaf dataset, the strokes in the estimated foreground content corresponding to Regime 1 (and 3) are relatively thicker when compared to those corresponding to Regime 2 across all the models (see Fig. 4(c)). Since the stroke widths are consistent with the ground truth, the PSNR continues to be higher for Regime 2. However, the F_{score} and PF_{score} of some models under Regime 2 and 3 are less than Regime 1. This is attributed to the presence of minor discontinuities in the strokes obtained under Regime 2 when compared to that of Regime 1 (see Fig. 4(c)). However, on average, the overall performance of each of the eight deep models is the highest when trained under Regime 2 across the three different test datasets, as evident from Table 2. Further, Fig. 4(c) depicts that models under Regime 1 and 3 fail to segment out the thread punched out through the document. However, almost every

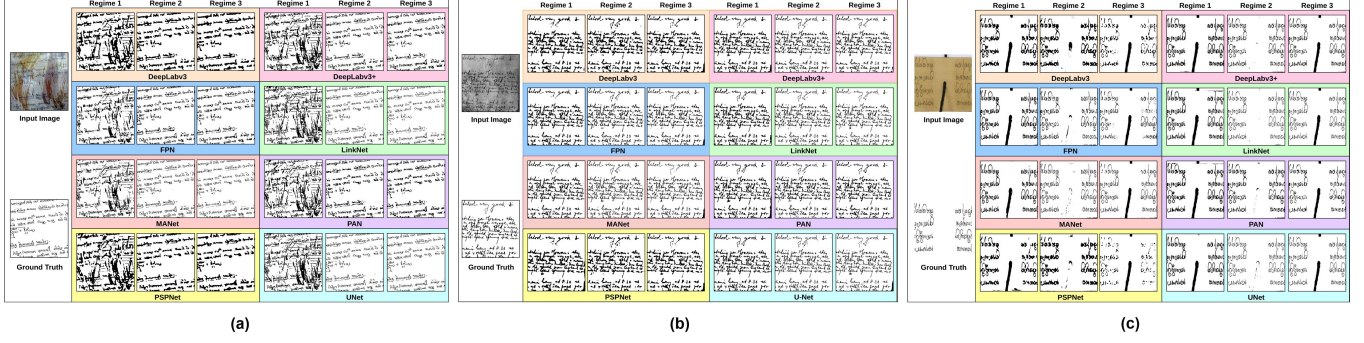


Fig. 4. Qualitative result on (a) the LS-HDIB test set (b) Bickley Diary dataset, and (c) Palm Leaf Manuscript dataset.

deep model trained on the LS-HDIB dataset (Regime 2) precisely identifies the appropriate foreground content even in the presence of such background artifacts.

For the DIBCO test set, the performance under Regime 1 is better than Regime 2 (Table 2). This is inline with the expectation as the networks are trained on DIBCO train set itself. Interestingly, Regime 3 offers the best performance across different models and test sets indicating that the LS-HDIB dataset when augmented with the standard DIBCO datasets enhances the network performance.

Metric	Dataset	DeepLabv3	DeepLabv3++	FPN	LinkNet	MANet	PANet	PSPNet	U-Net
F-score	Bickley Diary	0.64 0.66 0.60	0.71 0.77 0.75	0.72 0.77 0.74	0.71 0.82 0.80	0.72 0.84 0.82	0.72 0.74 0.69	0.58 0.65 0.60	0.74 0.83 0.77
	Palm Leaf	0.49 0.48 0.48	0.59 0.57 0.58	0.58 0.59 0.57	0.62 0.60 0.61	0.62 0.62 0.62	0.57 0.57 0.44	0.47 0.48 0.61	0.62 0.60 0.61
	DIBCO	0.48 0.43 0.52	0.68 0.57 0.62	0.60 0.46 0.61	0.66 0.58 0.72	0.68 0.60 0.71	0.56 0.41 0.58	0.50 0.43 0.46	0.67 0.58 0.72
	LS-HDIB	0.44 0.54 0.53	0.52 0.68 0.66	0.54 0.68 0.67	0.56 0.86 0.85	0.57 0.87 0.85	0.51 0.65 0.64	0.42 0.53 0.52	0.55 0.87 0.86
	Bickley Diary	0.64 0.66 0.60	0.76 0.79 0.77	0.75 0.79 0.75	0.80 0.87 0.82	0.79 0.89 0.80	0.71 0.76 0.73	0.59 0.66 0.61	0.78 0.89 0.82
	Palm Leaf	0.50 0.48 0.48	0.60 0.58 0.59	0.57 0.63 0.58	0.61 0.63 0.62	0.62 0.64 0.64	0.58 0.57 0.57	0.48 0.48 0.44	0.63 0.62 0.63
PF-score	DIBCO	0.49 0.41 0.52	0.68 0.52 0.62	0.60 0.51 0.62	0.66 0.58 0.73	0.68 0.51 0.72	0.56 0.41 0.59	0.50 0.43 0.47	0.67 0.58 0.72
	LS-HDIB	0.43 0.53 0.52	0.51 0.67 0.65	0.53 0.66 0.66	0.56 0.87 0.86	0.56 0.88 0.86	0.51 0.64 0.63	0.42 0.52 0.52	0.54 0.88 0.87
	Bickley Diary	10.07 10.06 9.69	12.32 12.75 12.40	12.11 12.69 12.29	13.16 14.42 13.41	12.83 14.82 13.49	11.55 12.10 11.64	9.61 10.01 9.78	12.89 14.64 13.24
	Palm Leaf	9.02 8.93 9.19	10.53 11.57 11.34	10.68 11.47 11.36	11.04 12.45 12.26	11.27 12.31 12.24	10.41 11.12 11.03	9.04 8.99 9.36	11.11 12.41 12.32
PSNR	DIBCO	8.26 6.75 8.38	9.76 8.97 10.47	9.63 8.74 10.34	10.84 10.96 12.54	11.13 9.54 12.36	9.01 10.53 9.74	7.90 6.56 8.39	11.06 11.01 12.56
	LS-HDIB	7.58 9.39 9.11	9.29 12.15 11.88	9.83 12.05 12.08	9.63 16.99 16.63	9.86 17.09 16.74	9.26 11.58 11.33	7.49 9.37 9.10	9.05 17.29 17.07

Table 2. F_{score} , PF_{score} , and PSNR evaluated over eight deep models under three different regimes: **Regime 1**, **Regime 2**, and **Regime 3** over test datasets.

Given that nearly all the models have performed the best when trained on LS-HDIB dataset across all the three test datasets, we further investigate the effect of dataset size on the network performance. For LS-HDIB test dataset, the F_{score} , PF_{score} , and PSNR are observed to increase with the dataset size, as shown in Fig. 5. This indicates the requirement of a large-scale dataset to span complex real-world scenarios encountered in the handwritten document images. Fig. 5 shows that the performance over the Bickley Diary and Palm Leaf dataset peaks at the dataset size of 20K and 5K, respectively. Overall, we have established that the deep models are more robust to various degradation effects and page styles encountered in the real world when trained on the LS-HDIB dataset and necessitate

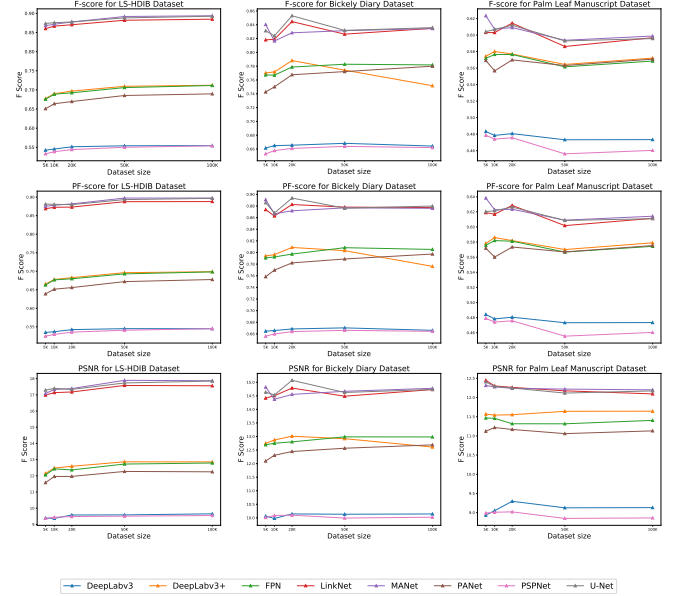


Fig. 5. Effect of varying dataset size on the model performance evaluated over the three test datasets.

the need for such a dataset.

Note. Owing to the space constraints, we provide more qualitative results, detailed dataset statistics, training and validation logs for different models (for better selectivity), and the accompanying code on our website¹.

5. CONCLUSION

We propose a large-scale dataset (LS-HDIB) for handwritten document image binarization and a simple yet effective method to generate it. When trained on the LS-HDIB dataset, different deep models can generalize better on unseen document images with a wide variety of degradations encountered in our day-to-day lives. This is possible due to the inherent diversity of the proposed dataset. Further, this work highlights that the fundamental image processing algorithms can be used as practical tools to support the existing deep-learning-based methods in producing significantly better results. In our case, we use adaptive thresholding and mixed gradient-based seamless cloning to generate this large-scale dataset.

¹<https://kaustubh-sadekar.github.io/LS-HDIB/>

6. REFERENCES

- [1] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 9, no. 1, pp. 62–66, 1979.
- [2] Pierre D Wellner, "Adaptive thresholding for the digitaldesk," *Xerox, EPC1993-110*, pp. 1–19, 1993.
- [3] Chris Tensmeyer and T. Martinez, "Document image binarization with fully convolutional neural networks," *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, vol. 01, pp. 99–104, 2017.
- [4] Quang Nhat Vo, Soo Hyung Kim, Hyung Jeong Yang, and Gueesang Lee, "Binarization of degraded document images based on hierarchical deep supervised network," *Pattern Recogn.*, vol. 74, no. C, pp. 568–586, Feb. 2018.
- [5] Sergey Levine, Peter Pastor, Alex Krizhevsky, Julian Ibarz, and Deirdre Quillen, "Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection," *The International Journal of Robotics Research*, vol. 37, no. 4-5, pp. 421–436, 2018.
- [6] Wen-Hsiang Tsai, "Moment-preserving thresholding: A new approach," *Computer vision, graphics, and image processing*, vol. 29, no. 3, pp. 377–393, 1985.
- [7] Basilios Gatos, Ioannis Pratikakis, and Stavros J Perantonis, "Adaptive degraded document image binarization," *Pattern recognition*, vol. 39, no. 3, pp. 317–327, 2006.
- [8] Joan Pastor-Pellicer, S España-Boquera, Francisco Zamora-Martínez, M Zeshan Afzal, and Maria Jose Castro-Bleda, "Insights on the use of convolutional neural networks for document image binarization," in *International Work-Conference on Artificial Neural Networks*. Springer, 2015, pp. 115–126.
- [9] Jorge Calvo-Zaragoza, Gabriel Viglienconi, and Ichiro Fujinaga, "Pixel-wise binarization of musical documents with convolutional neural networks," in *2017 Fifteenth IAPR International Conference on Machine Vision Applications (MVA)*. IEEE, 2017, pp. 362–365.
- [10] Muhammad Zeshan Afzal, Joan Pastor-Pellicer, Faisal Shafait, Thomas M Breuel, Andreas Dengel, and Marcus Liwicki, "Document image binarization using lstm: A sequence learning approach," in *Proceedings of the 3rd international workshop on historical document imaging and processing*, 2015, pp. 79–84.
- [11] Chris Tensmeyer, Mike Brodie, Daniel Saunders, and Tony Martinez, "Generating realistic binarization data with generative adversarial networks," in *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 2019, pp. 172–177.
- [12] Xujun Peng, Huaigu Cao, and Prem Natarajan, "Using convolutional encoder-decoder for document image binarization," in *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 2017, vol. 1, pp. 708–713.
- [13] Florian Westphal, Niklas Lavesson, and Håkan Grahn, "Document image binarization using recurrent neural networks," in *2018 13th IAPR International Workshop on Document Analysis Systems (DAS)*. IEEE, 2018, pp. 263–268.
- [14] A. K. Bhunia, A. K. Bhunia, A. Sain, and P. P. Roy, "Improving document binarization via adversarial noise-texture augmentation," in *2019 IEEE International Conference on Image Processing (ICIP)*, 2019, pp. 2721–2725.
- [15] B. Gatos, K. Ntirogiannis, and I. Pratikakis, "Icdar 2009 document image binarization contest (dibco 2009)," in *2009 10th International Conference on Document Analysis and Recognition*, 2009, pp. 1375–1382.
- [16] I. Pratikakis, B. Gatos, and K. Ntirogiannis, "H-dibco 2010 - handwritten document image binarization competition," in *2010 12th International Conference on Frontiers in Handwriting Recognition*, 2010, pp. 727–732.
- [17] I. Pratikakis, B. Gatos, and K. Ntirogiannis, "Icdar 2011 document image binarization contest (dibco 2011)," in *2011 International Conference on Document Analysis and Recognition*, 2011, pp. 1506–1510.
- [18] I. Pratikakis, B. Gatos, and K. Ntirogiannis, "Icfhr 2012 competition on handwritten document image binarization (h-dibco 2012)," in *2012 International Conference on Frontiers in Handwriting Recognition*, 2012, pp. 817–822.
- [19] I. Pratikakis, B. Gatos, and K. Ntirogiannis, "Icdar 2013 document image binarization contest (dibco 2013)," in *2013 12th International Conference on Document Analysis and Recognition*, 2013, pp. 1471–1476.
- [20] K. Ntirogiannis, B. Gatos, and I. Pratikakis, "Icfhr2014 competition on handwritten document image binarization (h-dibco 2014)," in *2014 14th International Conference on Frontiers in Handwriting Recognition*, 2014, pp. 809–813.
- [21] Ioannis Pratikakis, Konstantinos Zagoris, George Barlas, and Basilios Gatos, "Icfhr2016 handwritten document image binarization contest (h-dibco 2016)," 10 2016, pp. 619–623.
- [22] Ioannis Pratikakis, Konstantinos Zagoris, George Barlas, and Basilios Gatos, "Icdar2017 competition on document image binarization (dibco 2017)," 11 2017, pp. 1395–1403.
- [23] Ioannis Pratikakis, Konstantinos Zagori, Panagiotis Kaddas, and Basilios Gatos, "Icfhr 2018 competition on handwritten document image binarization (h-dibco 2018)," in *2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR)*, 2018, pp. 489–493.
- [24] H. Z. Nafchi, S. M. Ayatollahi, R. F. Moghaddam, and M. Cheriet, "An efficient ground truthing tool for binarization of historical manuscripts," in *2013 12th International Conference on Document Analysis and Recognition*, 2013, pp. 807–811.
- [25] Patrick Pérez, Michel Gangnet, and Andrew Blake, "Poisson image editing," in *ACM SIGGRAPH 2003 Papers*, New York, NY, USA, 2003, SIGGRAPH '03, p. 313–318, Association for Computing Machinery.
- [26] Liang-Chieh Chen, G. Papandreou, Florian Schroff, and Hartwig Adam, "Rethinking atrous convolution for semantic image segmentation," *ArXiv*, vol. abs/1706.05587, 2017.
- [27] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Computer Vision – ECCV 2018*, 2018.
- [28] Selim Seferbekov, Vladimir Iglovikov, Alexander Buslaev, and Alexey Shvets, "Feature pyramid network for multi-class land segmentation," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2018, pp. 272–2723.
- [29] Abhishek Chaurasia and Eugenio Culurciello, "Linknet: Exploiting encoder representations for efficient semantic segmentation," *2017 IEEE Visual Communications and Image Processing (VCIP)*, Dec 2017.
- [30] Tongle Fan, Guanglei Wang, Yan Li, and Hongrui Wang, "Ma-net: A multi-scale attention network for liver and tumor segmentation," *IEEE Access*, vol. 8, pp. 179656–179665, 2020.
- [31] Hanchao Li, Pengfei Xiong, Jie An, and Lingxue Wang, "Pyramid attention network for semantic segmentation," 2018.
- [32] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia, "Pyramid scene parsing network," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6230–6239.
- [33] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI 2015*. 2015, pp. 234–241, Springer International Publishing.
- [34] Pavel Golik, Patrick Doetsch, and Hermann Ney, "Cross-entropy vs. squared error training: a theoretical and experimental comparison," 08 2013, pp. 1756–1760.
- [35] A. F. Bickley, "Bickley diary dataset," 1926.
- [36] Jean, Mickaël, Setiawan, Kesiman, and MarcOgier, "Cfhr 2016 competition on the analysis of handwritten text in images of balinese palm leaf manuscripts," 10 2016, pp. 596–601.
- [37] K. Ntirogiannis, B. Gatos, and I. Pratikakis, "Performance evaluation methodology for historical document image binarization," *IEEE Transactions on Image Processing*, vol. 22, no. 2, pp. 595–609, 2013.