

DocDiff: Document Enhancement via Residual Diffusion Models

Zongyuan Yang
Beijing University of Posts and
Telecommunications
Beijing, China
yangzongyuan0@bupt.edu.cn

Baolin Liu
Beijing University of Posts and
Telecommunications
Beijing, China
baolin@bupt.edu.cn

Yongping Xiong*
Beijing University of Posts and
Telecommunications
Beijing, China
ypxiong@bupt.edu.cn

Lan Yi
Beijing University of Posts and
Telecommunications
Beijing, China

Guibin Wu
Beijing University of Posts and
Telecommunications
Beijing, China

Xiaojun Tang
Beijing University of Posts and
Telecommunications
Beijing, China

Ziqi Liu
Beijing University of Posts and
Telecommunications
Beijing, China

Junjie Zhou
Beijing University of Posts and
Telecommunications
Beijing, China

Xing Zhang
Artificial Intelligence Lab of China
Resources Digital Technology
Guangdong, China

ABSTRACT

Removing degradation from document images not only improves their visual quality and readability, but also enhances the performance of numerous automated document analysis and recognition tasks. However, existing regression-based methods optimized for pixel-level distortion reduction tend to suffer from significant loss of high-frequency information, leading to distorted and blurred text edges. To compensate for this major deficiency, we propose DocDiff, the first diffusion-based framework specifically designed for diverse challenging document enhancement problems, including document deblurring, denoising, and removal of watermarks and seals. DocDiff consists of two modules: the **Coarse Predictor (CP)**, which is responsible for recovering the primary low-frequency content, and the **High-Frequency Residual Refinement (HRR)** module, which adopts the diffusion models to predict the residual (high-frequency information, including text edges), between the ground-truth and the CP-predicted image. DocDiff is a compact and computationally efficient model that benefits from a well-designed network architecture, an optimized training loss objective, and a deterministic sampling process with short time steps. Extensive experiments demonstrate that DocDiff achieves state-of-the-art (SOTA) performance on multiple benchmark datasets, and can significantly enhance the readability and recognizability of degraded document images. Furthermore, our proposed HRR module in pre-trained DocDiff is plug-and-play and ready-to-use, with only 4.17M parameters. It greatly sharpens the text edges generated by SOTA

deblurring methods without additional joint training. **Available codes:** <https://github.com/Royalvice/DocDiff>.

CCS CONCEPTS

• **Computing methodologies** → **Computer vision**.

KEYWORDS

Document Enhancement; Conditional Diffusion Models; Frequency Separation; Document Analysis

ACM Reference Format:

Zongyuan Yang, Baolin Liu, Yongping Xiong, Lan Yi, Guibin Wu, Xiaojun Tang, Ziqi Liu, Junjie Zhou, and Xing Zhang. 2023. DocDiff: Document Enhancement via Residual Diffusion Models. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3581783.3611730>

1 INTRODUCTION

Document images are widely used in multi-media applications. The vulnerability to degradation is one of the challenges encountered in processing such highly structured data that differs significantly in pixel distribution from natural scene images. As an important pre-processing step, document enhancement aims to restore degraded document images to improve their readability and the performance of OCR systems [15]. In this paper, we focus on three major document enhancement tasks: document deblurring, document denoising and binarization, i.e., to remove fragmented noise from documents, such as smears and bleed-throughs, and watermark and seal removal. (Figure 6 in Appendix shows several degraded document images addressed in this paper.)

The major challenges for document enhancement are noise elimination and pixel-level text generation with low latency on high resolution document images. Specially, the presence of diverse noise types in document images, comprising of both global noise such as blurring and local noise such as smears, bleed-throughs, and seals, along with their potential combination, poses a significant challenge for noise elimination. Moreover, the task of generating

*Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '23, October 29–November 3, 2023, Ottawa, ON, Canada

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0108-5/23/10...\$15.00

<https://doi.org/10.1145/3581783.3611730>

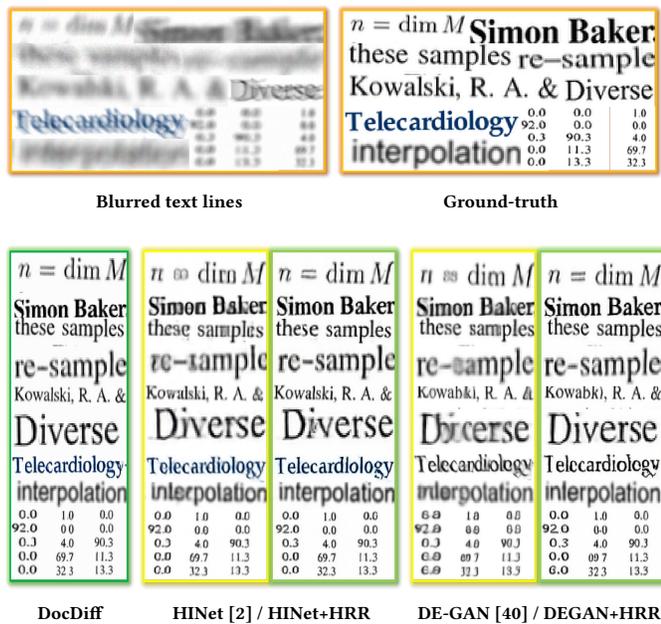


Figure 1: Compared to DE-GAN [40] and HINet [2], DocDiff generates sharper text edges. Our proposed HRR module effectively mitigates the problem of distorted and blurred characters generated by regression-based methods, regardless of whether they are designed for natural or document scenes. Note that the HRR module does not require additional joint training with DE-GAN [40] and HINet [2]. Its weights are derived from the pre-trained DocDiff.

text-laden images is fraught with difficulty. Unlike images that depict natural scenes, the high-frequency information of text-laden images is mostly concentrated on the text edges. The slightest erroneous pixel modifications at the text edges have the potential to alter the semantic meaning of a character, rendering it illegible or unrecognizable by OCR systems. Thus, document enhancement does not prioritize generation diversity, which differs from the pursuit of recovering multiple distinct denoised images in natural scenes [37, 43]. In practice, a typical document image contains millions of pixels. To ensure the efficiency of the entire document analysis system, pre-processing speed is crucial, which demands models to be as lightweight as possible.

Currently, existing document enhancement methods [39–41] are deep learning-based regression methods. Due to the problem of "regression to the mean", these methods [39–41] optimized with pixel-level loss produce blurry and distorted text edges. Additionally, due to the existence of numerous non-text regions in high-resolution document images, GAN-based methods [40] are prone to mode collapse when trained on local patches. In natural scenes, many diffusion-based methods [35, 43] try to restore degraded images with more details. However, there are challenges in directly applying these methods to document enhancement. Foremost, their high training costs and excessively long sampling steps make them

difficult to be practically implemented in document analysis systems. For these methods, shorter inference time implies the need for shorter sampling steps, which can lead to substantial performance degradation. Besides, the generation diversity of these methods can result in character inconsistency between the conditions and the sampled images.

Considering these problems, we transform document enhancement into the task of conditional image-to-image generation and propose DocDiff, which consists of two modules: the Coarse Predictor (CP) and the High-Frequency Residual Refinement (HRR) module. The CP takes degraded document images as input and approximately restores their clean versions. The HRR module leverages diffusion models to sharpen the text edges produced by CP accurately and efficiently. To avoid enhancing the generation quality at the expense of altering the characters and speed the inference, we adjust the optimization objective during training and adopt a short-step deterministic sampling strategy during inference. Specifically, we allow the HRR module to learn the distribution of residuals between the ground-truth images and the CP-predicted images (conditions). While ensuring consistency in the reverse diffusion process, we enable the HRR to directly predict the original residual rather than the added noise. This allows DocDiff to produce reasonable images that are highly correlated with the conditions in the first few steps of the reverse diffusion, based on the premise of using a channel-wise concatenation conditioning scheme [34, 35, 43]. While sacrificing generation diversity, which is not a critical factor for document enhancement, this operation considerably reduces the number of diffusion steps required for sampling and results in a reduced range of potential outputs generated by the conditional diffusion model, which effectively mitigates the production of distorted text edges and the replacement of characters. Overall, DocDiff undergoes end-to-end training with frequency separation through the joint use of pixel loss and modified diffusion model loss and applies the deterministic short-step sampling. DocDiff strikes the balance between innovation and practicality, achieving outstanding performance with a tiny and efficient network that contains only 8M parameters.

Experimental results on the three benchmark datasets (Document Deblurring [12] and (H-)DIBCO [29, 32]) demonstrate that DocDiff achieves SOTA performance in terms of perceptual quality for low-level deblurring task and competitive performance for high-level binarization task. More importantly, DocDiff achieves competitive deblurring performance with only 5 sampling steps. For the task of watermark and seal removal, we generated paired datasets using in-house document images. Experimental results demonstrate that DocDiff can effectively remove watermarks and seals while preserving the covered characters. Specifically, for seal removal, DocDiff trained on the synthesized dataset shows promising performance in real-world scenarios. Ablation experiments demonstrate the effectiveness of the HRR module on sharpening blurred characters generated by regression methods, as shown in Fig. 1.

In summary, the contributions of our paper are as follows:

- We present a novel framework, named DocDiff. To the best of our knowledge, it is the first diffusion-based method which

is specifically designed for diverse challenging document enhancement tasks.

- We propose a plug-and-play High-Frequency Residual Refinement (HRR) module to refine the generation of text edges. We demonstrate that the HRR module is capable of directly and substantially enhancing the perceptual quality of deblurred images generated by regression methods without requiring any additional training.
- DocDiff is a tiny, flexible, efficient and train-stable generative model. Our experiments show that DocDiff achieves competitive performance with only 5 sampling steps. Compared with non-diffusion-based methods [15, 39, 40, 51], DocDiff’s inference is fast with the same level of performance. Additionally, DocDiff is trained inexpensively, avoids mode collapse and can enhance both handwritten and printed document images at any resolution.
- Adequate ablation studies and comparative experiments show that DocDiff achieves competitive performance on the tasks of deblurring, denoising, watermark and seal removal on documents. Our results highlight the benefits of various components of DocDiff, which collectively contribute to its superior performance.

2 RELATED WORKS

2.1 Document Enhancement

The pixel distribution in document images differs significantly from that of natural scene images, and most document images have a resolution greater than 1k. Therefore, it is crucial to develop specialized enhancement models for document scenarios to handle the degradation of different types of documents efficiently and robustly. The currently popular document enhancement methods [15, 19, 39, 40] are predominantly based on deep learning regression methods. These methods aim to achieve higher PSNR by minimizing \mathcal{L}_1 or \mathcal{L}_2 pixel loss. However, distortion metrics like PSNR only partially align with human perception [1]. This problem is particularly noticeable in document scenarios (see details depicted in Fig. 10 in Appendix). Although GAN-based methods [40, 41, 51] utilize a combination of content and adversarial losses to generate images with sharp edges, training GANs on high-resolution document datasets can be challenging because the cropped patches typically have a significant number of identical patterns, which increases the risk of mode collapse.

2.2 Diffusion-based Image-to-image

Recently, Diffusion Probabilistic Models (DPMs) [10, 38] have been widely adopted for conditional image generation [18, 21, 33–35, 37, 43?]. Saharia et al. [35] present SR3, which adapts diffusion models to image super-resolution. Saharia et al. [34] propose Palette, which is a multi-task image-to-image diffusion model. Palette has demonstrated the excellent performance of diffusion models in the field of conditional image generation, including colorization, inpainting and JPEG restoration. In [18, 37, 43], they all utilize the prediction-refinement framework where the diffusion models are used to predict the residual. Different from [18, 21, 37], DocDiff is trained end-to-end.

Although effective for denoising natural images, they are not specifically designed for document images. Methods [18, 21, 33–35, 37] cannot handle images of arbitrary resolution. Due to their large networks, long sampling steps and large number of cropped patches, these lead to prohibitively long inference time on high resolution document images. Although Method [43] is capable of processing images of any resolution and has a relatively small network structure, its training method of predicting noise and stochastic sampling strategy produce diverse characters and distorted text edges, which is not suitable for document enhancement.

3 METHODOLOGY

The overall architecture of DocDiff is shown in Fig. 2. DocDiff consists of two modules: the Coarse Predictor (CP) and the High-Frequency Residual Refinement (HRR) module. Due to the fixed pattern of document images and their varying resolutions, we adopt a compact U-Net structure for the CP and HRR module, modified from [10]. We replace the self-attention layer in the middle with four layers of dilated convolutions to increase the receptive field, and remove the remaining self-attention and normalization layers. To reduce computational complexity, we compress the parameters of the CP and HRR module to 4.03M and 4.17M respectively, while ensuring performance, making them much smaller than existing document enhancement methods [39, 40, 45, 51]. Overall, DocDiff is a fully convolutional model that employs a combination of pixel loss and diffusion model loss, facilitating end-to-end training with frequency separation.

3.1 Coarse Predictor

The objective of the Coarse Predictor C_θ is to approximately restore a degraded document image y into its clean version x_{gt} at the pixel level. The $\mathcal{L}_{\text{pixel}}$ can be defined as the mean square error between the coarse prediction x^C and the x_{gt} :

$$x^C = C_\theta(y) \quad (1)$$

$$\mathcal{L}_{\text{pixel}} = \mathbb{E} \|x^C - x_{gt}\|_2 \quad (2)$$

During text pixel generation, the Coarse Predictor can effectively restore the primary content of the text, but it may not be able to accurately capture the high-frequency information in the text edges. This leads to significant blurring at the edges of the text. As shown in Fig. 1, this is a well-known limitation of CNN-based regression methods due to the problem of "regression to the mean". It is challenging to address this issue by simply cascading more convolutional layers.

3.2 High-Frequency Residual Refinement Module

To address the problem mentioned above, we introduce the High-Frequency Residual Refinement (HRR) module, which is capable of generating samples from the learned posterior distribution. The core component of HRR module is a Denoiser f_θ which leverages DPMs to estimate the distribution of residuals between the ground-truth images and the images generated by the CP. In contrast to prior research [18, 21, 37], we design the HRR module to address not only the "regression to the mean" flaw found in a single regression model (in this case, the CP), but also to be effective across a

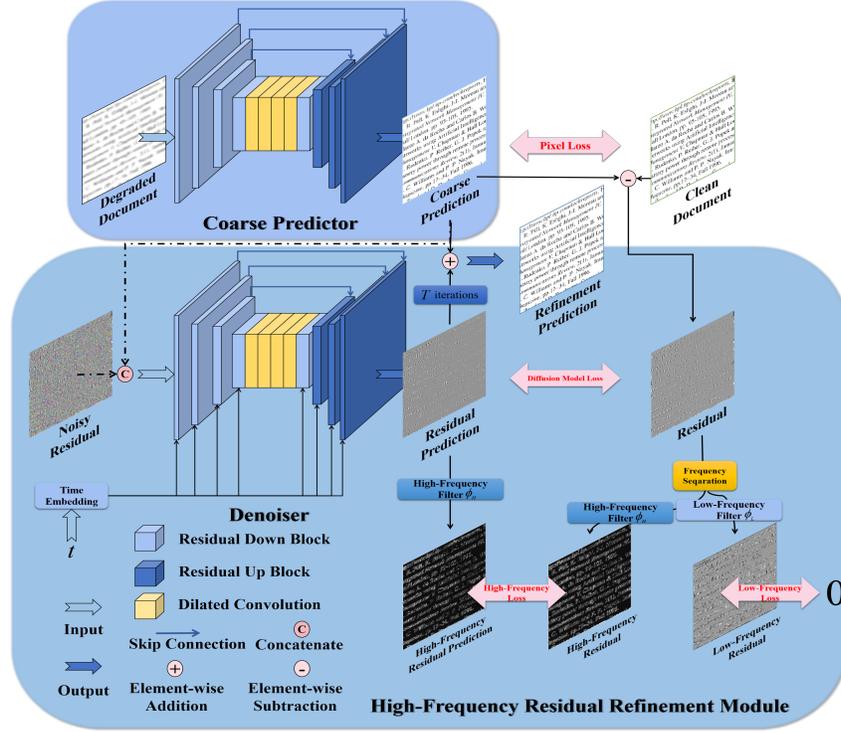


Figure 2: A overview of the proposed DocDiff for document enhancement. Take document deblurring as an example.

variety of regression methods. To this end, We perform end-to-end joint training of both CP and HRR modules, rather than separately. By this way, the HRR module can dynamically adjust and capture more patterns. Extensive experiments show that this training strategy effectively enhances the sharpness of characters generated by different regression-based deblurring methods [2, 17, 39, 40, 48], without requiring joint training like [21, 37].

3.2.1 *Denoiser* f_{θ} . Followed by [10, 38], the HRR module executes the **forward noise-adding process** and the **reverse denoising process** to model the residual distributions.

Forward noise-adding process: Given the clean document image x_{gt} and its approximate estimate x^C , we calculate their residuals x_{res} . We assign x_0 as x_{res} , and then sequentially introduce Gaussian noise based on the time step, as follows:

$$x_0 = x_{res} = x_{gt} - x^C \quad (3)$$

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t} x_{t-1}, (1 - \alpha_t) \mathbf{I}) \quad (4)$$

$$q(x_{1:T} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1}) \quad (5)$$

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I}) \quad (6)$$

where α_t is a hyperparameter between 0 and 1 that controls the variance of the added Gaussian noise at each time step, for all $t = 1, \dots, T$, and $\alpha_0 = 1, \bar{\alpha}_t = \prod_{i=0}^t \alpha_i$. There are no learnable

parameters in the forward process, and $x_{1:T}$ have the same size as x_0 . With the reparameterization trick, x_t can be written as:

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (7)$$

Reverse denoising process: The reverse process transforms Gaussian noise back into the residual distributions conditioned on x^C . We can write the reverse diffusion step:

$$q(x_{t-1} | x_t, x_0) = \mathcal{N}(x_{t-1}; \mu_t(x_t, x_0), \beta_t(x_t, x_0) \mathbf{I}) \quad (8)$$

where $\mu_t(x_t, x_0)$ and $\beta_t(x_t, x_0)$ are the mean and variance, respectively. Followed by [38], we perform the deterministic reverse process $q(x_{t-1} | x_t, x_0)$ with zero variance and the mean can be computed:

$$\mu_t(x_t, x_0) = \sqrt{\bar{\alpha}_{t-1}} x_0 + \sqrt{1 - \bar{\alpha}_{t-1}} \cdot \underbrace{\frac{x_t - \sqrt{\bar{\alpha}_t} x_0}{\sqrt{1 - \bar{\alpha}_t}}}_{\epsilon} \quad (9)$$

$$\beta_t(x_t, x_0) = 0 \quad (10)$$

Given the Denoiser f_{θ} , the posterior $q(x_{t-1} | x_t, x_0)$ can be parameterized:

$$p_{\theta}(x_{t-1} | x_t, x^C) = q(x_{t-1} | x_t, f_{\theta}(x_t, t, x^C)) \quad (11)$$

We emphasize the significance of the condition x^C in the conditional distribution $p_{\theta}(x_{t-1} | x_t, x^C)$. At each time step, it is crucial to sample the residual that closely relate to x^C at the pixel level of the characters.

The Denoiser f_θ can be trained to predict the original data x_0 or the added noise ϵ . To increase the diversity of generated natural images, existing methods [10, 34, 35, 43] typically predict the added noise ϵ . The prediction of x_0 and ϵ is equivalent in unconditional generation, as they can be transformed into each other through Eq.7. However, for conditional generation, predicting x_0 and predicting ϵ are not equivalent under the premise of using a channel-wise concatenation conditioning scheme [34, 35, 43] to introduce the condition x^C . When predicting ϵ , the denoiser can only learn from the noisy x_t . But when predicting x_0 , the denoiser can also learn from the conditional channels (x^C). This sacrifices diversity but significantly improves the generation quality of the first few steps in the reverse process. To this end, we train the Denoiser f_θ to directly predict x_0 , which aligns with the goals of document enhancement. The training objective is to minimize the distance between $p_\theta(x_{t-1} | x_t, x^C)$ and the true posterior $q(x_{t-1} | x_t, x_0)$:

$$\mathcal{L}_{DM} = \mathbb{E} \left\| x_0 - f_\theta \left(\sqrt{\alpha_t} x_{\text{res}} + \sqrt{1 - \alpha_t} \epsilon, t, \hat{x}^C \right) \right\|_2 \quad (12)$$

where \hat{x}^C is a clone of x^C in memory and does not participate in gradient calculations. Exactly, the gradient from the loss only flows through x_{res} from f_θ to C_θ .

Given the trained C_θ and f_θ , integrating the above equation we finally have the deterministic reverse process:

$$\hat{x}_{\text{res}} = f_\theta(x_t, t, C_\theta(y)) \quad (13)$$

$$x_{t-1} = \sqrt{\alpha_{t-1}} \hat{x}_{\text{res}} + \sqrt{1 - \alpha_{t-1}} \cdot \frac{x_t - \sqrt{\alpha_t} \hat{x}_{\text{res}}}{\sqrt{1 - \alpha_t}} \quad (14)$$

3.2.2 Frequency Separation Training. Numerous prior studies [5, 6, 37] have demonstrated that processing high and low frequency information separately enhances the quality and level of detail in generated images of natural scenes. To further refine the generation quality, DocDiff is trained with frequency separation. Specifically, we employ simple linear filters to separate the residuals into the low and high frequencies. As the spatial-domain addition is equivalent to the frequency-domain addition, frequency separation can be written directly as:

$$x = \phi_L * x + \phi_H * x \quad (15)$$

where ϕ_L is the low-pass filter and ϕ_H is the high-pass filter. In practice, to extract residuals along the text edges, we commonly set ϕ_H as the Laplacian kernel. The low-frequency information is obtained according to Eq.15. Our approach differs from [37, 44] in that it does not require performing the Fast Fourier Transform (FFT) and parameterizing frequency separation in the frequency domain. The high-frequency information in document images is primarily concentrated at the text edges. Leveraging this prior knowledge by using the Laplacian kernel as a high-pass filter not only simplifies training time consumption but also proves to be highly effective.

Our **goal** is to maximize the capacity of the Denoiser f_θ to restore the missing high-frequency information in the Coarse Predictor C_θ prediction, while minimizing the task burden of f_θ through the support of C_θ . In this perspective, both C_θ and f_θ necessitate the restoration of distinct high and low frequencies, yet they exhibit dissimilar specializations. Specifically, C_θ predominantly reconstructs low-frequency information, while f_θ specializes in restoring

high-frequency details:

$$\mathcal{L}^{\text{low}} = \mathbb{E} \|\phi_L * (x^C - x_{\text{gt}})\|_2 \quad (16)$$

$$\mathcal{L}^{\text{high}} = \mathbb{E} \left\| \phi_H * \left(x_0 - f_\theta(x_t, t, \hat{x}^C) \right) \right\|_2 \quad (17)$$

Thus, the combined loss for C_θ and f_θ , together with the overall loss for DocDiff, are given by:

$$\mathcal{L}_{\text{pixel}}^{\text{low}} = \mathcal{L}_{\text{pixel}} + \beta_0 \mathcal{L}^{\text{low}} \quad (18)$$

$$\mathcal{L}_{DM}^{\text{high}} = \mathcal{L}_{DM} + \beta_0 \mathcal{L}^{\text{high}} \quad (19)$$

$$\mathcal{L}_{\text{total}} = \beta_1 \mathcal{L}_{\text{pixel}}^{\text{low}} + \mathcal{L}_{DM}^{\text{high}} \quad (20)$$

Algorithm 1 outlines the complete training and inference processes of DocDiff.

Algorithm 1 DocDiff

Require : C_θ : Coarse Predictor, f_θ : Denoiser,
 y : Degraded document image, ϕ_H : High frequency filter,
 x_{gt} : Clean document image, $\alpha_0:T$: Noise schedule.

Training :

- 1: **repeat**
- 2: $(y, x_{\text{gt}}) \sim q(y, x_{\text{gt}})$, $t \sim \text{Uniform}\{1, \dots, T\}$, $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 3: $x^C = C_\theta(y)$ ▷ Coarse prediction
- 4: $x_t = \sqrt{\alpha_t}(x_{\text{gt}} - x^C) + \sqrt{1 - \alpha_t}\epsilon$ ▷ Forward diffusion
- 5: $\hat{x}_{\text{res}} = f_\theta(x_t, t, x^C)$ ▷ Residual prediction
- 6: Take a gradient descent step on
 $\mathcal{L}_{\text{total}}(x^C, \hat{x}_{\text{res}}, x_{\text{gt}}, \phi_H; C_\theta, f_\theta)$
 ▷ Frequency separation training; see Eq. 20

7: **until** converged

Sampling:

- 1: $x^C = C_\theta(y)$ ▷ Coarse prediction
 - 2: $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 3: **for** $t = T, \dots, 1$ **do**
 - 4: $\hat{x}_{\text{res}} = f_\theta(x_t, t, x^C)$
 - 5: $x_{t-1} = \sqrt{\alpha_{t-1}} \hat{x}_{\text{res}} + \frac{\sqrt{1 - \alpha_{t-1}}(x_t - \sqrt{\alpha_t} \hat{x}_{\text{res}})}{\sqrt{1 - \alpha_t}}$
 ▷ Deterministic reverse diffusion
 - 6: **end for**
 - 7: **return** $x^C + x_0$ ▷ Return high-frequency residual refinement
-

4 EXPERIMENTS AND RESULTS

4.1 Datasets and Implementation details

To address various document enhancement tasks, we train and evaluate our models on distinct datasets. For a fair comparison, we use the source codes provided by the authors and follow the same experimental settings.

Document Deblurring: We train and evaluate DocDiff on the widely-used Document Deblurring Dataset [12], which includes 66,000 pairs of clean and blurry 300×300 patches extracted from diverse pages of different documents. Each blur kernel is distinct. We randomly select 30,000 patches for training and 10,000 patches for testing.

Document Denoising and Binarization: We evaluate DocDiff on two of the most challenging datasets from the annual (Hand-written) Document Image Binarization Competition ((H-)DIBCO) [7, 23, 25–32]: H-DIBCO'18 [29] and DIBCO'19 [32]. Followed by [41, 45, 47], the training set includes the remaining years of DIBCO

Table 1: Quantitative ablation study results on the Document Deblurring Dataset [12]. Best values are highlighted in red, second best are highlighted in blue. (CP: Coarse Predictor, CR: Cascade Refinement Module, HRR: High-Frequency Residual Refinement Module, EMA: Exponential Moving Average, FS: Frequency Separation Training.)

Structure components		Training techniques		Method				Perceptual				Distortion			
CP	CR	HRR	EMA	FS	Resolutions	Sampling		Parameters	NR		FR				
					Non-native	Native	Stochastic	Deterministic	MANIQA \uparrow	MUSIQ \uparrow	DISTS \downarrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	
✓			✓			✓	-	-	4.03M	0.6525	46.15	0.0951	0.0766	24.66	0.9574
✓	✓		✓	✓		✓	-	-	8.06M	0.6584	45.39	0.0688	0.0824	24.74	0.9610
✓		✓			✓			✓	8.20M	0.6900	50.16	0.0671	0.0492	20.98	0.9025
✓		✓	✓		✓			✓	8.20M	0.6917	50.21	0.0648	0.0499	20.43	0.8998
✓		✓	✓	✓	✓		✓		8.20M	0.6706	50.05	0.1778	0.1481	18.72	0.8507
✓		✓	✓	✓	✓			✓	8.20M	0.6971	50.31	0.0636	0.0474	20.46	0.9006
✓	✓	✓	✓	✓		✓		✓	8.20M	0.7174	50.62	0.0611	0.0307	23.28	0.9505

datasets [7, 23, 25–28, 30, 31], the Noisy Office Database [49], the Bickley Diary dataset [3], Persian Heritage Image Binarization Dataset (PHIDB) [20], and the Synchromedia Multispectral dataset (S-MS) [9].

Watermark and Seal Removal: Dense watermarks and variable seals greatly affect the readability and recognizability of covered characters. There is limited research on this issue in the document analysis community and a paucity of publicly available benchmark datasets. Thus, we synthesized paired datasets (document image with dense watermarks and seals and its corresponding clean version, see synthetic details in Synthetic Datasets section and Fig. 9 in Appendix) using in-house data for training and testing. For the discussion of experimental results on the synthetic datasets, please refer to Section C in Appendix.

We jointly train the CP and HRR modules by minimizing the \mathcal{L}_{total} in Eq. 20 with $\beta_0 = 2$ and $\beta_1 = 0.5$. The total time steps T are set to 100. We use random 128×128 crops with a batchsize of 64 during training, and evaluate DocDiff on both non-native (larger size patches or full-size images) and native (128×128 crop-predict-merge) resolutions and different sampling steps (5, 20, 50 and 100). For data augmentation, we perform random rotation and horizontal flip. The number of training iterations is one million.

4.2 Evaluation Metrics

We employ the SOTA no-reference (NR) image quality assessment (IQA) methods, including MUSIQ [14] that is sensitive to varying sizes and aspect ratios and MANIQA [46] that won the NTIRE 2022 NR-IQA Challenge [8], as well as widely-used full-reference (FR) IQA methods including LPIPS [50] and DISTS [4], to evaluate the reconstruction quality of document images. We still compute PSNR and SSIM [42] for completeness, although not the primary metrics.

For high-level binarization tasks, we evaluate methods on three evaluation metrics commonly used in competitions, including the FMeasure (FM), pseudo-FMeasure (p-FM) [22], and PSNR. For removal tasks, we evaluate methods on four metrics: MANIQA [46], LPIPS [50], PSNR and SSIM [42].

4.3 Document Deblurring

4.3.1 Ablation Study. We conduct the ablation experiments on the Document Deblurring Dataset [12] to verify the benefits of

proposed components of DocDiff: High-Frequency Residual Refinement Module (HRR), Frequency Separation Training (FS). Additionally, we investigate the impact of resolutions (native or non-native), sampling method (stochastic or deterministic), and predicted target (x_0 or ϵ) on the performance of the model. Table 1 shows the quantitative results. (see qualitative results in Fig. 11 in Appendix.)

HRR: The HRR module effectively improves the perceptual quality by sharpening text edges. To further prove that the effectiveness of HRR is not solely due to an increase in parameters, we cascade an identical Unet structure (CR) behind CP to transform the model into a two-stage regression method. Experimental results show that while simply cascading more encoder-decoder layers improve PSNR and SSIM, the perceptual quality remain poor with blurred text edges. This reiterates the effectiveness of the HRR module.

FS: As shown in Table 1, training with frequency separation demonstrates an improvement in the perceptual quality and a reduction in distortion, thus achieving a better Perception-Distortion trade-off. This decoupled strategy can effectively enhance the capacity of HRR module to recover high-frequency information.

Native or Non-native ? : Performing the crop-predict-merge strategy at native resolution can significantly reduce the distortion. However, inferring at non-native resolution (full-size image) also yields competitive perceptual quality and distortion. In practice, this is a time-quality trade-off. For instance, at 300×300 resolution, using 128×128 crop-predict-merge inference requires processing 60% more ineffective pixels compared to full-size inference.

Predict x_0 or ϵ , Stochastic or Deterministic ? : We employ the HRR module to predict the ϵ and perform original stochastic sampling [10]. Except for NR-IQA, experimental results show inferior performance in FR-IQA, PSNR, and SSIM for this approach. On one hand, the short sampling step (100) restricts the capability of the noise prediction-based models. On the other hand, stochastic sampling results can lead to poor integration of generated characters with the given conditions, which leads to increased occurrence of substitution characters. This is also the reason for the high NQ-IQA score and low FQ-IQA score. These issues are addressed by combining the prediction of x_0 and deterministic sampling. Moreover, we notice that the approach based on predicting x_0 can effectively generate high-quality images with 5 sampling steps (see Tab. 3 and Tab. 7 in Appendix), whereas the approach relying on predicting ϵ fails to achieve the same level of performance.

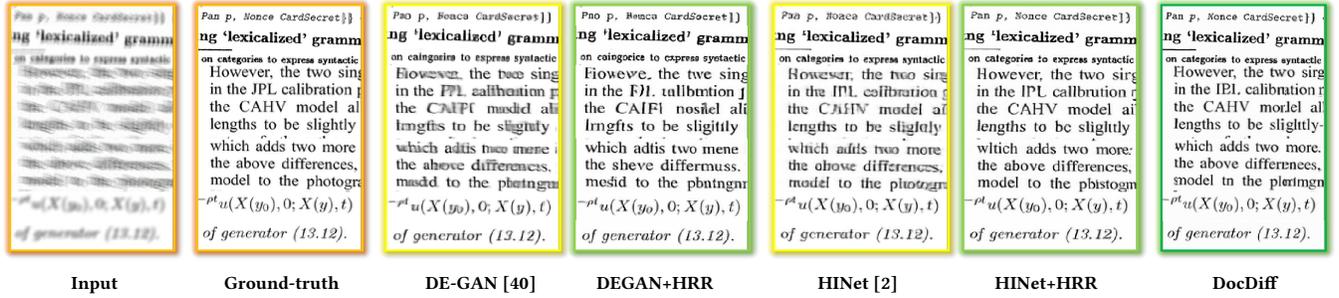


Figure 3: Qualitative representative results of text line deblurring on the Document Deblurring Dataset [12].

Table 2: Results on the Document Deblurring Dataset [12]. The bold numbers represent the improvement of the HRR module over the regression-based baseline. Note that the weights of the HRR module come from DocDiff and no joint training with the baseline has been performed.

Method	Perceptual				Distortion	
	NR		FR		PSNR \uparrow	SSIM \uparrow
	MANIQA \uparrow	MUSIQ \uparrow	DISTS \downarrow	LPIPS \downarrow		
Natural scenes						
DeBlurGAN-v2 [17], ICCV2019	0.6967	50.25	0.1014	0.0991	20.76	0.8731
+HRR	0.7097	50.66	0.0952	0.0989	20.62	0.8736
Pec. of Better than GT	38.78%	48.83%				
Pages Better than DeBlurGAN-v2	68.53%	82.76%				
MPRNet[48], CVPR2021	0.6675	47.52	0.1555	0.0900	21.27	0.8803
+HRR	0.6852	49.87	0.1384	0.0887	20.86	0.8768
Pec. of Better than GT	17.87%	40.12%				
Pages Better than MPRNet	57.04%	87.04%				
HINet [2], CVPR2021	0.6836	47.59	0.1232	0.1163	24.15	0.9164
+HRR	0.7041	50.44	0.0963	0.0987	23.44	0.9158
Pec. of Better than GT	31.07%	46.67%				
Pages Better than HINet	56.91%	91.82%				
Document scenes						
DE-GAN [40], TPAMI2020	0.6546	46.75	0.0968	0.0843	22.30	0.9155
+HRR	0.6973	50.36	0.0776	0.0696	21.21	0.9114
Pec. of Better than GT	23.48%	45.34%				
Pages Better than DE-GAN	79.89%	95.40%				
DocEnTr [39], ICPR2022	0.5821	46.53	0.1802	0.2225	22.66	0.9130
+HRR	0.6637	51.84	0.1378	0.1653	21.65	0.9142
Pec. of Better than GT	5.57%	68.57%				
Pages Better than DocEnTr	94.40%	93.71%				
GT	0.7207	51.03	0.0	0.0	∞	1.0

4.3.2 *Performance.* For a comprehensive comparison, we compare with SOTA document deblurring methods [39, 40] as well as natural scene deblurring methods [2, 17, 43, 48]. Table 3 shows quantitative results and Figure 3 shows qualitative results. DocDiff(Native) achieves the **best** MANIQA, DISTS, LPIPS, and SSIM metrics, while also achieving competitive PSNR and MUSIQ. Notably, we obtain the LPIPS of 0.0307, a **66% reduction** compared to MPRNet [48] and a **64% reduction** compared to DE-GAN [40]. DocDiff uses only **one-fourth** of the parameters used by those two methods. We also compare DocDiff with SOTA diffusion-based deblurring method [43], which predicts ϵ and samples stochastically. [43] can produce high-quality images (higher MUSIQ). However, its FR-IQA and distortion metrics are significantly worse than DocDiff’s at the same sampling step. Moreover, DocDiff (Non-native) outperforms MPRNet[48], HINet[2], and two document-scene methods [39, 40] in **all perceptual metrics** with only **5-step** sampling, while maintaining competitive distortion metrics.

Table 3: Document deblurring results on the Document Deblurring Dataset [12]. DocDiff outperforms state-of-the-art deblurring regression methods for both natural and document scenes on all perceptual metrics, even at non-native resolutions, while maintaining competitive PSNR and SSIM scores. DocDiff-n means applying n-step sampling (T).

Method	Parameters	Perceptual				Distortion	
		NR		FR		PSNR \uparrow	SSIM \uparrow
		MANIQA \uparrow	MUSIQ \uparrow	DISTS \downarrow	LPIPS \downarrow		
Natural scenes							
DeBlurGAN-v2 [17], ICCV2019	67M	0.6967	50.25	0.1014	0.0991	20.76	0.8731
MPRNet[48], CVPR2021	34M	0.6675	47.52	0.1555	0.0900	21.27	0.8803
HINet [2], CVPR2021	86M	0.6836	47.59	0.1232	0.1163	24.15	0.9164
Whang et al. [43], CVPR2022	33M	0.6898	50.86	0.0830	0.0750	19.89	0.8742
Document scenes							
DE-GAN [40], TPAMI2020	31M	0.6546	46.75	0.0968	0.0843	22.30	0.9155
DocEnTr [39], ICPR2022	67M	0.5821	46.53	0.1802	0.2225	22.66	0.9130
DocDiff (Non-native)-5	8.20M	0.6873	47.92	0.0907	0.0582	22.17	0.9223
DocDiff (Non-native)-100	8.20M	0.6971	50.31	0.0636	0.0474	20.46	0.9006
DocDiff (Native)-100	8.20M	0.7174	50.62	0.0611	0.0307	23.28	0.9505
GT	-	0.7207	51.03	0.0	0.0	∞	1.0

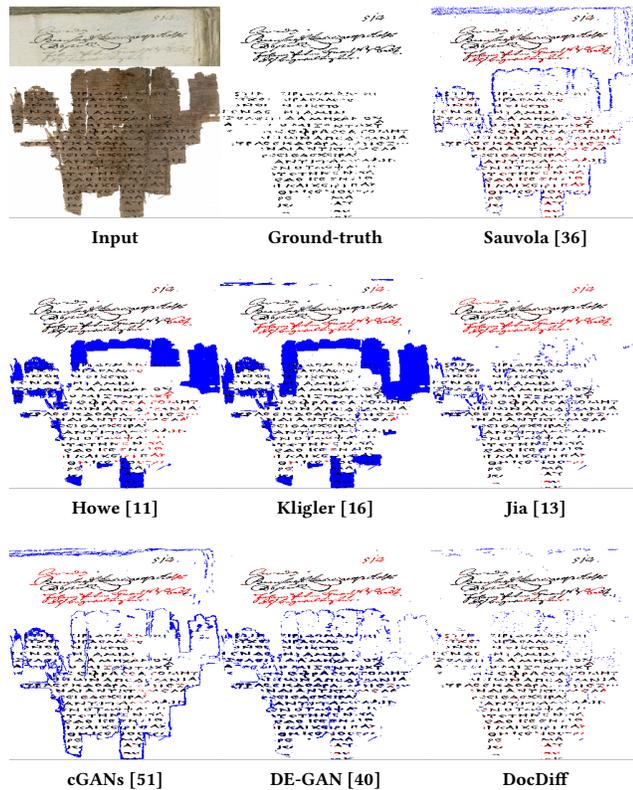
To validate the universality of the HRR module, we refine the output of baselines [2, 17, 39, 40, 48] directly using the pre-trained HRR module. As shown in Tab. 2, after refinement, **all perceptual metrics** are improved. Specifically, the DISTS of HINet [2] decrease by **22%**, and the MANIQA of DocEnTr [39] increase by **14%**. We calculate the percentage of samples in which NR-IQA show better performance compared to the GT and baselines after refinement. On average, in terms of MUSIQ, **90%** of the samples show improvement over baselines, and **50%** over GT. As shown in Figs. 1 and 3, DocDiff provides the most clear and accurate capability of character pixel restoration. After refined by the HRR module, the text edges of the baselines become sharp, but wrong characters still exist.

4.4 Document Denoising and Binarization

We compare with threshold-based methods [11, 13, 16, 24, 36] and SOTA methods [40, 45, 51]. Quantitative and qualitative results are shown in Tab. 4 and Fig. 4, respectively. On the H-DIBCO’18, our method do not achieve SOTA performance. However, our F-Measure is 10.52% higher than that of DE-GAN [40], and we also have competitive performance in terms of p-FM and PSNR. Due to the absence of papyrus material in training sets and the presence of a large amount of newly introduced noise in DIBCO’19, this dataset poses a great challenge. DocDiff outperforms existing methods on DIBCO’19 with an F-Measure improvement of 5.75% over D²

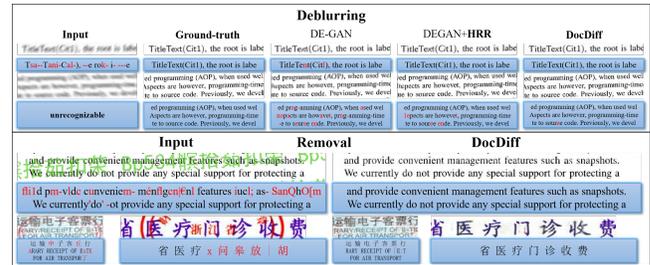
Table 4: Results of document binarization on H-DIBCO 2018 [29] and DIBCO 2019 [32]. Best values are bold.

Method	Parameters	H-DIBCO'18			DIBCO'19		
		FM \uparrow	p-FM \uparrow	PSNR \uparrow	FM \uparrow	p-FM \uparrow	PSNR \uparrow
Otsu [24]	-	51.45	53.05	9.74	47.83	45.59	9.08
Sauvola [36], PR2000	-	67.81	74.08	13.78	51.73	55.15	13.72
Howe [11], IJDAR2013	-	80.84	82.85	16.67	48.20	48.37	11.38
Jia et al. [13], PR2018	-	76.05	80.36	16.90	55.87	56.28	11.34
Kligler et al. [16], CVPR2018	-	66.84	68.32	15.99	53.49	53.34	11.23
1st rank of contest	-	88.34	90.24	19.11	72.88	72.15	14.48
cGANs [51], PR2019	103M	87.73	90.60	18.37	62.33	62.89	12.43
DE-GAN [40], TPAMI2020	31M	77.59	85.74	16.16	55.98	53.44	12.29
D ² BFormer [45], IF2023	194M	88.84	93.42	18.91	67.63	66.69	15.05
DocDiff	8.20M	88.11	90.43	17.92	73.38	75.12	15.14

**Figure 4: Binarization results of two example images on H-DIBCO 2018 [29] and DIBCO 2019 [32]. Text pixels classified as background are highlighted in red, whereas background pixels classified as text are highlighted in blue.****Table 5: Average runtimes (seconds/megapixel) of different methods.**

Method	MPRNet [48]	DE-GAN [40]	DocDiff(Non-native)-5	DocDiff(Non-native)-100
Runtime	0.57	0.82	0.33	5.69

BFormer [45] and 11.05% over cGANs [51]. Remarkably, this is achieved using significantly fewer parameters, with only 1/23 and 1/12 the parameter count of D² BFormer [45] and cGANs [51], respectively.

**Figure 5: Qualitative results for Tesseract recognition of various text lines. The red characters indicate recognition errors or missed characters. Best viewed with zoom-in.**

Diffusion models are notorious for their time complexity of inference. DocDiff (Non-native)-5 achieve competitive performance in removing blur and watermarks. We compare the time complexity of our methods with MPRNet [48] and DE-GAN [40] under the same hardware environment. However, due to different frameworks being used (Tensorflow for DE-GAN [40] while PyTorch for MPRNet [48] and our methods), the speed comparison is only approximate. The results are shown in Tab. 5. DocDiff (Non-native)-5 offers great efficiency and performance, making it an ideal choice for various document image enhancement tasks.

4.5 OCR Evaluation

we compare the OCR performance on degraded and enhanced documents using a set of 50 text patches. This set includes 30 blurred patches from the Document Deblurring Dataset [12], 10 patches with watermarks, and 10 patches with seals. We use Tesseract OCR to recognize those patches. Highly blurred images are barely recognizable. After applying DE-GAN [40], DE-GAN[40]+HRR, and DocDiff for deblurring, the character error rates of the enhanced versions are reduced to 13.7%, 7.6%, and 4.4%, respectively. For the removal task, the character error rates of the original images and the ones enhanced with DocDiff are 28.7% and 1.8%, respectively. The recognition results on some patches are shown in Fig. 5.

5 CONCLUSIONS

In this paper, we propose a novel and unified framework, DocDiff, for various document enhancement tasks. DocDiff significantly improves the perceptual quality of reconstructed document images by utilizing a residual prediction-based conditional diffusion model. For the deblurring task, our proposed HRR module is ready-to-use, which effectively sharpens the text edges generated by regression methods [2, 17, 39, 40, 48] to enhance the readability and recognizability of the text. Compared to non-diffusion-based methods, DocDiff achieves competitive performance with only 5 steps of sampling, and is lightweight, which greatly optimizes its inference time complexity. We believe that DocDiff establishes a strong benchmark for future work.

REFERENCES

- [1] Yochai Blau and Tomer Michaeli. 2018. The perception-distortion tradeoff. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 6228–6237.

- [2] Liangyu Chen, Xin Lu, Jie Zhang, Xiaojie Chu, and Chengpeng Chen. 2021. Hinet: Half instance normalization network for image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 182–192.
- [3] Fanbo Deng, Zheng Wu, Zheng Lu, and Michael S. Brown. 2010. BinarizationShop: A User-Assisted Software Suite for Converting Old Documents to Black-and-White. In *Proceedings of the 10th Annual Joint Conference on Digital Libraries - JCDL '10*. ACM Press, Gold Coast, Queensland, Australia, 255.
- [4] Keyan Ding, Kede Ma, Shiqi Wang, and Eero P. Simoncelli. 2022. Image Quality Assessment: Unifying Structure and Texture Similarity. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 5 (2022), 2567–2581. <https://doi.org/10.1109/TPAMI.2020.3045810>
- [5] Manuel Fritsche, Shuhang Gu, and Radu Timofte. 2019. Frequency separation for real-world super-resolution. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. IEEE, 3599–3608.
- [6] Dario Fuoli, Luc Van Gool, and Radu Timofte. 2021. Fourier space losses for efficient perceptual image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2360–2369.
- [7] Basilis Gatos, Konstantinos Ntirogiannis, and Ioannis Pratikakis. 2009. ICDAR 2009 Document Image Binarization Contest (DIBCO 2009). In *2009 10th International Conference on Document Analysis and Recognition*. IEEE, Barcelona, Spain, 1375–1382.
- [8] Jinjin Gu, Haoming Cai, and Chao et al. Dong. 2022. NTIRE 2022 Challenge on Perceptual Image Quality Assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*. 951–967.
- [9] Rachid Hedjam, Hossein Ziaei Nafchi, Reza Farrahi Moghaddam, Margaret Kalaska, and Mohamed Cheriet. 2015. ICDAR 2015 Contest on MultiSpectral Text Extraction (MS-TEX 2015). In *2015 13th International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, Tunis, Tunisia, 1181–1185.
- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.
- [11] Nicholas R. Howe. 2013. Document Binarization with Automatic Parameter Tuning. *International Journal on Document Analysis and Recognition (IJ DAR)* 16, 3 (Sept. 2013), 247–258.
- [12] Michal Hradš, Jan Kotera, Pavel Zemeč, and Filip Šroubek. 2015. Convolutional neural networks for direct text deblurring. In *Proceedings of BMVC*, Vol. 10.
- [13] Fuxi Jia, Cunzhaoh Shi, Kun He, Chunheng Wang, and Baihua Xiao. 2018. Degraded Document Image Binarization Using Structural Symmetry of Strokes. *Pattern Recognition* 74 (Feb. 2018), 225–240.
- [14] Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. 2021. MUSIQ: Multi-Scale Image Quality Transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 5148–5157.
- [15] Sana Khamekhem Jemni, Mohamed Ali Souibgui, Yousri Kessentini, and Alicia Fornés. 2022. Enhance to Read Better: A Multi-Task Adversarial Network for Handwritten Document Image Enhancement. *Pattern Recognition* 123 (March 2022), 108370.
- [16] Netanel Kligler, Sagi Katz, and Ayellet Tal. 2018. Document enhancement using visibility detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2374–2382.
- [17] Orest Kupyń, Tetiana Martyniuk, Junru Wu, and Zhangyang Wang. 2019. Deblurgan-v2: Deblurring (orders-of-magnitude) faster and better. In *Proceedings of the IEEE/CVF international conference on computer vision*. 8878–8887.
- [18] Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting Chen. 2022. Srdiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing* 479 (2022), 47–59.
- [19] Yun-Hsuan Lin, Wen-Chin Chen, and Yung-Yu Chuang. 2020. Bedsr-net: A deep shadow removal network from a single document image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12905–12914.
- [20] Hossein Ziaei Nafchi, Seyed Morteza Ayatollahi, Reza Farrahi Moghaddam, and Mohamed Cheriet. 2013. An Efficient Ground Truthing Tool for Binarization of Historical Manuscripts. In *2013 12th International Conference on Document Analysis and Recognition*. IEEE, Washington, DC, USA, 807–811.
- [21] Axi Niu, Kang Zhang, Trung X Pham, Jinqiu Sun, Yu Zhu, In So Kweon, and Yanning Zhang. 2023. CDPMSR: Conditional Diffusion Probabilistic Models for Single Image Super-Resolution. *arXiv preprint arXiv:2302.12831* (2023).
- [22] K. Ntirogiannis, B. Gatos, and I. Pratikakis. 2013. Performance Evaluation Methodology for Historical Document Image Binarization. *IEEE Transactions on Image Processing* 22, 2 (Feb. 2013), 595–609.
- [23] Konstantinos Ntirogiannis, Basilis Gatos, and Ioannis Pratikakis. 2014. ICFHR2014 Competition on Handwritten Document Image Binarization (H-DIBCO 2014). In *2014 14th International Conference on Frontiers in Handwriting Recognition*. IEEE, Greece, 809–813.
- [24] Nobuyuki Otsu. 1979. A Threshold Selection Method from Gray-Level Histograms. *IEEE Transactions on Systems, Man, and Cybernetics* 9, 1 (1979), 62–66.
- [25] Ioannis Pratikakis, Basilis Gatos, and Konstantinos Ntirogiannis. 2010. H-DIBCO 2010 - Handwritten Document Image Binarization Competition. In *2010 12th International Conference on Frontiers in Handwriting Recognition*. IEEE, Kolkata, India, 727–732.
- [26] Ioannis Pratikakis, Basilis Gatos, and Konstantinos Ntirogiannis. 2011. ICDAR 2011 Document Image Binarization Contest (DIBCO 2011). In *2011 International Conference on Document Analysis and Recognition*. IEEE, Beijing, China, 1506–1510.
- [27] Ioannis Pratikakis, Basilis Gatos, and Konstantinos Ntirogiannis. 2012. ICFHR 2012 Competition on Handwritten Document Image Binarization (H-DIBCO 2012). In *2012 International Conference on Frontiers in Handwriting Recognition*. IEEE, Bari, Italy, 817–822.
- [28] Ioannis Pratikakis, Basilis Gatos, and Konstantinos Ntirogiannis. 2013. ICDAR 2013 Document Image Binarization Contest (DIBCO 2013). In *2013 12th International Conference on Document Analysis and Recognition*. IEEE, Washington, DC, USA, 1471–1476.
- [29] Ioannis Pratikakis, Konstantinos Zagori, Panagiotis Kaddas, and Basilis Gatos. 2018. ICFHR 2018 Competition on Handwritten Document Image Binarization (H-DIBCO 2018). In *2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR)*. IEEE, Niagara Falls, NY, 489–493.
- [30] Ioannis Pratikakis, Konstantinos Zagoris, George Barlas, and Basilis Gatos. 2016. ICFHR2016 Handwritten Document Image Binarization Contest (H-DIBCO 2016). In *2016 15th International Conference on Frontiers in Handwriting Recognition (ICFHR)*. IEEE, Shenzhen, China, 619–623.
- [31] Ioannis Pratikakis, Konstantinos Zagoris, George Barlas, and Basilis Gatos. 2017. ICDAR2017 Competition on Document Image Binarization (DIBCO 2017). In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, Kyoto, 1395–1403.
- [32] Ioannis Pratikakis, Konstantinos Zagoris, Xenofon Karagiannis, Lazaros Tsochatzidis, Tanmoy Mondal, and Isabelle Marthot-Santaniello. 2019. ICDAR 2019 Competition on Document Image Binarization (DIBCO 2019). In *2019 International Conference on Document Analysis and Recognition (ICDAR)*. 1547–1556.
- [33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10684–10695.
- [34] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. 2022. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*. 1–10.
- [35] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J. Fleet, and Mohammad Norouzi. 2023. Image Super-Resolution via Iterative Refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 4 (2023), 4713–4726. <https://doi.org/10.1109/TPAMI.2022.3204461>
- [36] J. Sauvola and M. Pietikäinen. 2000. Adaptive Document Image Binarization. *Pattern Recognition* 33, 2 (Feb. 2000), 225–236.
- [37] Shuyao Shang, Zhengyang Shan, Guangxing Liu, and Jinglin Zhang. 2023. ResDiff: Combining CNN and Diffusion Model for Image Super-Resolution. *arXiv preprint arXiv:2303.08714* (2023).
- [38] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502* (2020).
- [39] Mohamed Ali Souibgui, Sanket Biswas, Sana Khamekhem Jemni, Yousri Kessentini, Alicia Fornés, Josep Lladós, and Umapada Pal. 2022. DocEnTr: an end-to-end document image enhancement transformer. In *2022 26th International Conference on Pattern Recognition (ICPR)*. IEEE, 1699–1705.
- [40] Mohamed Ali Souibgui and Yousri Kessentini. 2020. De-gan: A conditional generative adversarial network for document enhancement. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 3 (2020), 1180–1191.
- [41] Sungsho Suh, Jihun Kim, Paul Lukowicz, and Yong Oh Lee. 2022. Two-Stage Generative Adversarial Networks for Binarization of Color Document Images. *Pattern Recognition* 130 (Oct. 2022), 108810.
- [42] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing* 13, 4 (2004), 600–612. <https://doi.org/10.1109/TIP.2003.819861>
- [43] Jay Whang, Mauricio Delbracio, Hossein Talebi, Chitwan Saharia, Alexandros G Dimakis, and Peyman Milanfar. 2022. Deblurring via stochastic refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16293–16303.
- [44] Junde Wu, Huihui Fang, Yu Zhang, Yehui Yang, and Yanwu Xu. 2022. MedSegDiff: Medical Image Segmentation with Diffusion Probabilistic Model. *arXiv preprint arXiv:2211.00611* (2022).
- [45] Mingming Yang and Songhua Xu. 2023. A Novel Degraded Document Binarization Model through Vision Transformer Network. *Information Fusion* 93 (2023), 159–173.
- [46] Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan Gong, Mingdeng Cao, Jiahao Wang, and Yujiu Yang. 2022. MANIQA: Multi-Dimension Attention Network for No-Reference Image Quality Assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*. 1191–1200.
- [47] Zongyuan Yang, Yongping Xiong, and Guibin Wu. 2023. GDB: Gated convolutions-based Document Binarization. *arXiv preprint arXiv:2302.02073* (2023).
- [48] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. 2021. Multi-stage progressive image

- restoration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 14821–14831.
- [49] F. Zamora-Martinez, S. España-Boquera, and M. J. Castro-Bleda. 2007. Behaviour-Based Clustering of Neural Networks Applied to Document Enhancement. In *Proceedings of the 9th International Work Conference on Artificial Neural Networks (IWANN'07)*. Springer-Verlag, Berlin, Heidelberg, 144–151.
- [50] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [51] Jinyuan Zhao, Cunzhao Shi, Fuxi Jia, Yanna Wang, and Baihua Xiao. 2019. Document Image Binarization with Cascaded Generators of Conditional Generative Adversarial Networks. *Pattern Recognition* 96 (Dec. 2019), 106968.

A DEGRADED DOCUMENTS

Figure 6 illustrates the three types of degraded documents that our work aims to enhance: documents with fragmented noise, blurry documents, and documents with seals and dense watermarks. Regarding the removal of seals, we mainly focus on red seals in the context of Chinese documents.



Figure 6: Examples of degraded document images.

B SYNTHETIC DATASETS

Figure 9 shows some examples from our synthetic datasets. The synthesized dense watermarks feature randomized text (including both Chinese and English characters and numbers), font, size, color, spacing, position, and angle. The opacity of the watermarks is randomly sampled between 0.7 to 0.95. We utilized our unified seal-segmentation method to extract the mask of the seals in real Chinese document scenes. These seals mostly come from our internal documents, with a small portion coming from the ICDAR 2023 Competition on Reading the Seal Title. Afterwards, we fused the seal masks into the background images in the same manner. In our developed datasets, documents covered with watermarks have a resolution of 1754×1240 (with 3000 512×512 patches allocated for training and 100 full images for testing), while those covered with seals have a resolution of 512×512 (with 3000 for training and 500 used for testing). Note that the watermarks, seals, and background images in the training sets and testing sets are independent.

Table 6: Watermark and seal removal results on our developed datasets.

Method	Watermark Removal				Seal Removal			
	Perceptual		Distortion		Perceptual		Distortion	
	MANQA↑	LPIPS↓	PSNR↑	SSIM↑	MANQA↑	LPIPS↓	PSNR↑	SSIM↑
MPRNet [48]	0.5253	0.0806	33.36	0.9497	0.5133	0.0913	32.76	0.9397
DE-GAN [40]	0.5190	0.1167	24.11	0.9106	0.4742	0.1115	19.62	0.7084
DocDiff(Non-native)-5	0.5263	0.0680	30.91	0.9637	0.5018	0.1162	31.07	0.9292
DocDiff(Non-native)-100	0.5267	0.0577	30.36	0.9576	0.5031	0.1152	30.67	0.9225
Ground Truth	0.5367	0.0	∞	1.0	0.5525	0.0	∞	1.0

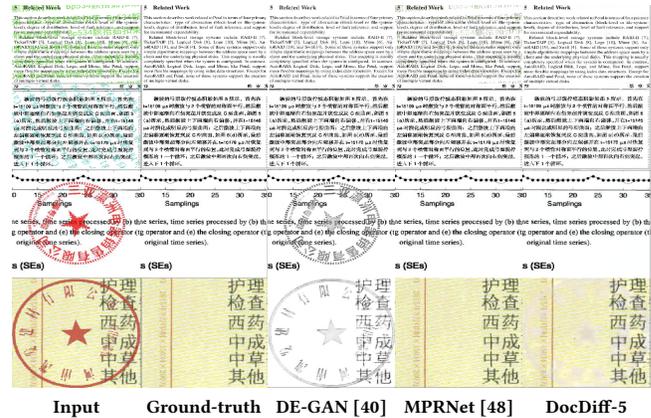


Figure 7: Qualitative watermark and seal removal results on synthetic datasets. While DocDiff is effective at removing watermarks and seals, MPRNet [48] displays better performance in restoring background on invoices.



Figure 8: Qualitative seal removal results in real Chinese invoice and document scenarios. Our proposed DocDiff model exhibits superior generalization ability and is resilient to noise.

C WATERMARK AND SEAL REMOVAL

Quantitative and qualitative results are shown in Tab. 6 and Fig. 7, respectively. For watermark removal, DocDiff (Non-native)-5 exhibits competitive performance as compared to MPRNet [48]. For

seal removal, MPRNet [48] perform better on synthetic datasets due to its well-designed feature extraction module that can effectively restore diverse invoice backgrounds. However, in real-world Chinese invoice and document scenarios, DocDiff demonstrate better generalization ability, as shown in Fig. 8. DE-GAN [40] is designed to take grayscale images as input and output, hence its performance on multi-color removal is bad.

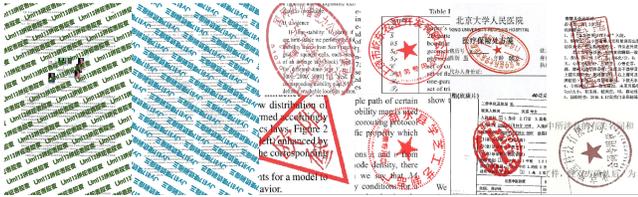


Figure 9: Examples of document images synthesized with dense watermarks and real seals.

D METRICS DESCRIPTION

As shown in Fig. 10, DE-GAN generates images with higher PSNR and SSIM scores but the character pixel-level edges are notably blurred, causing difficulty in reading for humans and recognition for OCR systems. With the enhancement of the HRR module, the character edges become much sharper and easier to recognize. The improvement in multiple perceptual metrics aligns with human perceptual quality.

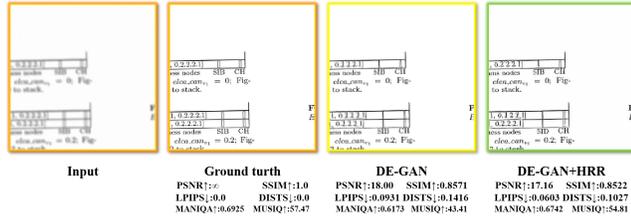


Figure 10: Comparison of different metrics. Best viewed with zoom-in.

E MORE DETAILS IN ABLATION STUDY

Figure 11 shows qualitative results in ablation study. We can get the following useful conclusions:

- Adding more encoder-decoder layers to optimize pixel loss in a cascade manner may not necessarily improve edge sharpening and legibility.
- Training with frequency separation can better restore text edges.
- Predicting added noise ϵ and applying short-step stochastic sampling can result in noisy sampled images with less sharp text edges.
- Inference at native resolution provides the best performance.

F DISCUSSION ABOUT SAMPLING STEPS

Followed by DDIM, the time steps during training are set to 100, while different sampling steps are used during inference including 5, 10, 20, 50, and 100. Table 7 shows quantitative results and Figure 12 shows qualitative results. As the sampling step increases, there is an observable trend where the perceptual quality of the image improves, but this comes at the cost of increased distortion. As shown in Fig. 12, the edge of the word "Expression" become distinguishable after 5 sampling steps, while the edge of the word "Therefore" only become clear after 50 sampling steps. While DocDiff can correctly restore a majority of characters, errors can still occur such as "Therefore" becoming "Therefor". We introduce a possible solution to this problem in the next section. Even with 100 sampling steps during inference is usually considered a low sampling step for general diffusion models. We emphasize again that DocDiff is able to restore relatively sharp text edges within 20 steps thanks to its training strategy of predicting original data x_0 and its deterministic sampling strategy.

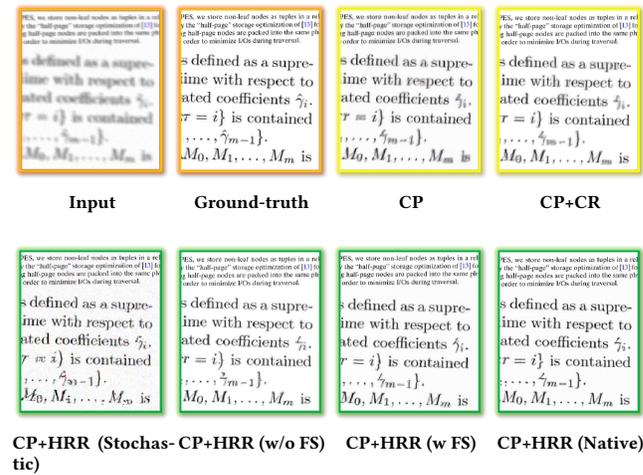


Figure 11: Qualitative ablation study results on the Document Deblurring Dataset. (CP: Coarse Predictor, CR: Cascade Refinement Module, HRR: High-Frequency Residual Refinement Module, FS: Frequency Separation Training.

Table 7: Quantitative results for different sampling steps on the Document Deblurring Dataset. DocDiff-n means applying n-step sampling (T).

	Perceptual				Distortion	
	MANQA \uparrow	MUSIQ \uparrow	DISTS \downarrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow
DocDiff (Non-native)-5	0.6873	47.92	0.0907	0.0582	22.17	0.9223
DocDiff (Non-native)-10	0.6878	47.93	0.0886	0.0592	22.16	0.9217
DocDiff (Non-native)-20	0.6890	47.99	0.0875	0.0608	22.13	0.9206
DocDiff (Non-native)-50	0.6912	48.13	0.0776	0.0565	21.88	0.9180
DocDiff (Non-native)-100	0.6971	50.31	0.0636	0.0474	20.46	0.9006

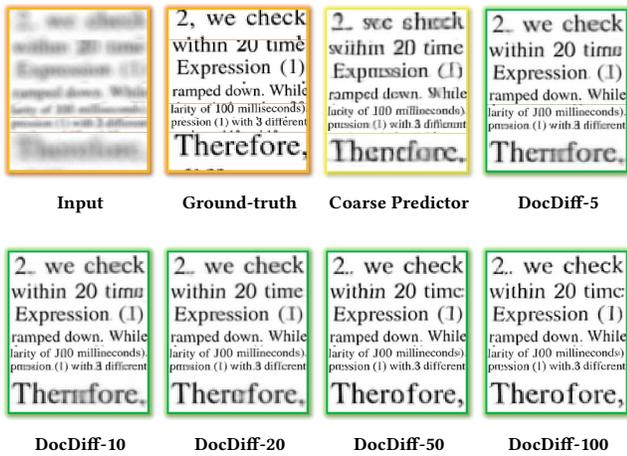


Figure 12: Qualitative sampling results for different time steps on the Document Deblurring Dataset.

G FUTURE WORKS

There are several potential solutions to overcome the limitations of our work. During the experiment, we observe that although DocDiff is able to generate sharp text edges in most cases, there are still instances where characters appear erroneous or distorted. To address this issue, one approach is to utilize the text prior in order to incorporate additional semantic information. Currently, there is a scarcity of paired large-scale document enhancement benchmark datasets in real-world scenarios. This leads to a lack of generalizability in practical applications. One approach to address this issue is to utilize DocDiff for assisted annotation. For instance, in the case of seal removal, humans can annotate on the results generated by DocDiff, thereby substantially reducing labelling complexity. Through multiple rounds of iteration and fine-tuning, the dataset can be expanded and the performance of the model can be improved.