# Restoration and Text Extraction for Enhanced Analysis of Vintage and Damaged Documents Using Deep Learning

Anand Magar, Rutuja Desai, Siddhi Deshmukh, Samarth Deshpande and Sakshi Dhamne

*Department of Computer Engineering, Vishwakarma Institute of Technology, Pune, India*

Abstract:     Restoration and analysis of vintage and damaged documents are crucial in preserving valuable historical and cultural records. In this paper, we present a novel approach that combines Document Enhancement Generative Adversarial Network (DE-GAN) for image restoration with Optical Character Recognition (OCR) for efficient text extraction. Our methodology focuses on restoring degraded documents by enhancing visual quality and legibility, allowing for more accurate text retrieval. By employing DE-GAN, we can mitigate various forms of document degradation, such as stains, faded ink and physical damage, while the OCR system effectively extracts the underlying text for further analysis. The proposed framework enhances the quality of historical document preservation and facilitates better data retrieval, ensuring that critical information is not lost due to document deterioration. Experimental results demonstrate significant improvements in both restoration quality and text extraction accuracy, making our method a robust solution for historical document analysis.

## 1 INTRODUCTION

This Historical and archival documents are an invaluable part of our collective heritage, but over time, many of these documents suffer from deterioration caused by environmental factors, physical damage, or improper storage. Such degradation often results in the loss of crucial information and poses a significant challenge to researchers and archivists who aim to preserve these documents for future generations. Traditional restoration techniques, while effective, can be labor-intensive and time-consuming. Moreover, manual text extraction from these degraded documents is prone to errors and inconsistencies.

Recent advances in deeplearning and computer vision have opened up new possibilities for automating the restoration and text extraction process. Generative Adversarial Networks (GANs), in particular, have shown great promise in various image restoration tasks due to their ability to generate high-quality images from degraded inputs. In this study, we introduce a DE-GAN-based approach to restore damaged documents, enhancing both their visual clarity and structural integrity. Paired with OCR technology, this enables the accurate extraction of text from the restored images, even in cases where traditional OCR techniques would struggle due to document quality.

The combination of DE-GAN for restoration and OCR for text extraction offers a comprehensive solution to the challenges posed by document degradation. Our approach not only helps preserve the visual quality of historical documents but also facilitates easier access to the text they contain, enabling researchers and archivists to perform more in-depth analysis on restored documents.

## 2 LITERATURE SURVEY

Literature survey has been carried out to understand innovative approaches.

Researchers have used deep learning method image inpainting and then valued performance of the model to identify how an algorithm can restore an image using many performance metrics. utilization of available power. EMS is crucial for optimal power balance in hybrid PV/Wind turbine systems. The objective of EMS is for improving transients, MPPT, EMS for grid. The flow chart for an EMS is shown in

Figure 4.. Authors have done a survey of the applications of deep learning in the field of image restoration and gives some ways for future research. (Shubekova, Beibitkyzy, et al. , 2023). Authors have used deep convolutional neural networks for research. Research include four restoration tasks include image inpainting, pixel interpolation, image deblurring and image denoising and used three different datasets. Results demonstrate that approach performs well on all models. (Gao and Grauman, 2017). In another study, two neural network structures were used to estimate background light and scene depth to enhance underwater image restoration. Experimental results on both synthetic and real underwater images highlight the effectiveness of the method. (Cao, Peng, et al. , 2018). Researchers developed a convolutional neural network (CNN) model to de-haze individual images, facilitating further restoration and improvement. The model's performance was evaluated using images and features from distinct regions, demonstrating its generalization capabilities. (Ramkumar, Ayyadurai, et al. , 2021). The development of speckle reduction techniques has advanced alongside progress in image restoration, with deep neural networks offering new solutions. These networks have surpassed earlier methods such as patches, sparsity, wavelet transforms, and total variation minimization in terms of restoration performance. (Denis, Dalsasso, et al. , 2021). A proposed method for digital mural restoration blocks high-priority sample blocks, prevents processing of areas with numerous unknown pixels, and minimizes error accumulation, achieving over 20% accuracy improvement compared to baseline methods. This technique excels in restoring the main structure and texture details of complex images with significant information loss. (Xiao, Zheng, et al. , 2023). A generative adversarial network (GAN) with embedded channels and spatial attention was introduced, increasing the perceptual field through expanded convolution. This allows the network to better capture image details and broken edges, leading to enhanced restoration performance, as shown by experimental results on public datasets. (Meng, Yu, et al. , 2022). A comparative study of methods for restoring noisy images, including the Wiener filter, wavelet methods, and Wiener filtering with BM3D techniques, was also conducted. The study analyzed the effectiveness of different color models like RGB and YCbCr in image segmentation, alongside CNN-based image classification, considering factors such as activation functions and pooling methods. (Vispute, Rajeswari, et al. , 2023). Lastly, a CNN regression model was employed to

learn enhancement parameters for different types of underwater images. Quantitative and qualitative experiments conducted on known datasets demonstrated impressive accuracy rates, achieving better results on PSNR and SSIM quality metrics compared to other techniques. (Martinho, Calvalcanti, et al. , 2024). Researchers in the field of handwritten text recognition (HTR) have focused on improving accuracy using deep learning models like LSTM. In one study, an LSTM-based HTR model, integrated with an OCR system, outperformed alternative methods on the IAM handwritten dataset, with the 2DLSTM approach showing superior results. (Nikitha, Geetha, et al. , 2020). A study on text recognition from real-time images compared two feature extraction methods: Histogram of Oriented Gradients (HOG) and Local Binary Patterns (LBP). After preprocessing and text localization, a CNN-based model was used for text recognition. LBP proved more effective in accuracy, precision, and F1 score compared to HOG. (Harsha, Kumar, et al. , 2022). Text recognition from natural images remains challenging due to complex backgrounds and irregular text arrangements. A deep learning model (DL-TRI) was developed to handle curved and perspective fonts, addressing these complexities more effectively than previous methods. (Shrivastava, Amudha, et al. , 2019). Handwriting recognition systems, which previously relied on optical character recognition (OCR), now utilize CNNs, LSTMs, and Connectionist Temporal Classification (CTC) for better accuracy. This approach was tested on the MNIST dataset using OpenCV and TensorFlow for image processing and word recognition, yielding improved performance. (Ansari, Kaur, et al. , 2022). Researchers are recognizing handwritten text, a two-phase method using CNNs was proposed. The first phase identifies the input, while the second phase classifies the language. The system, tested on the MNIST database, achieved 99% testing accuracy and 99.6% training accuracy. (Thilagavathy, Suresh, et al. , 2023). The architecture, design, and testing of a Handwritten Character Recognition System are presented, demonstrating the efficiency of neural networks in recognizing handwritten characters. The system converts images of handwritten notes from students and instructors into digital text for further use. (MS and A. D. R, 2023). Deep learning and machine learning techniques have shown significant promise in optical character recognition (OCR). A comprehensive overview of four key architectures—Support Vector Machines (SVM), Artificial Neural Networks (ANN), Naive Bayes, and Convolutional Neural Networks (CNN)—is provided, highlighting

their roles in improving OCR. (Sharma, Kaushik, et al. , 2020). The evolution of OCR due to advancements in deep learning algorithms is discussed, focusing on the introduction of Convolutional Recurrent Neural Networks (CRNN) combined with attention mechanisms to enhance text recognition performance. (N. J and S. K. A M , 2024). The increasing availability of GPUs and cloud platforms such as Google Cloud and Amazon Web Services has expanded computational power, facilitating the training of complex neural networks. This study details the design of an image segmentation-based handwritten character recognition system using OpenCV for image processing and TensorFlow for neural network training, implemented in Python. (Vaidya, Trivedi, et al. , 2018). Handwritten Text Recognition (HTR) has garnered significant attention due to its broad applications, yet it remains a challenging research area. The variability in personal handwriting and the characteristics of handwritten characters across different languages pose difficulties for HTR systems. Traditional approaches rely on Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) with the Connectionist Temporal Classification (CTC) objective function. However, attention-based sequence-to-sequence (Seq2Seq) models have emerged as a more flexible solution, particularly for the temporal nature of text. These models utilize attention mechanisms to focus on the most relevant features of the input, offering enhanced adaptability. This paper provides a comprehensive comparison of deep learning approaches for HTR and discusses the challenges that limit their effectiveness. (Teslya, Mohammed, et al. , 2022).

# 3 METHODOLOGY

The methodology of this research paper focuses on the use of Document Enhancement Generative Adversarial Network (DE-GAN) for restoring vintage and damaged documents, coupled with Optical Character Recognition (OCR) to extract text for analysis. The methodology is divided into the following steps: dataset preparation, GAN architecture, training process, and evaluation.

## 3.1 Dataset Preparation

We use a set of degraded document images from the DIBCO datasets (2009-2017) for training and

evaluation. The dataset consists of two types of images.

- **Degraded images (A):** These are the input images affected by various forms of damage such as stains, fading, and noise.

- **Clean images (B):** These images serve as the ground truth for the restoration process.

Before feeding the images into the model, the degraded and clean images are resized to a uniform size of $256 \times 256 \times 1$ to ensure compatibility with the GAN network architecture. Preprocessing is applied to convert the images to grayscale to simplify the restoration task.

## 3.2 GAN Architecture

The restoration process is carried out using a **Generative Adversarial Network (GAN)**, which consists of a **generator** and a **discriminator**:

- **Generator Model:** The generator's task is to restore degraded document images to their original, clean state. It takes as input a degraded image and outputs a restored version. The generator is built using a convolutional neural network with multiple layers, gradually increasing the resolution of the generated images. The largest layer in the generator has 1024 units to handle complex restorations.

- **Discriminator Model:** The discriminator's job is to distinguish between real clean images and those generated by the generator. It learns to classify the restored images as either real (clean) or fake (restored). This binary classification helps the generator improve over time. The discriminator is also built using a convolutional architecture and is periodically trained on both the real and generated images.

The generator and discriminator models are combined into a **GAN network**, where the generator tries to fool the discriminator by generating realistic images, and the discriminator attempts to classify them correctly.

## 3.3 Training Process

The GAN training is performed as follows:

**Adversarial Training:** The generator and discriminator are trained in alternating steps. In each epoch, the following procedure is followed:

- **Load Data:** For each image in the dataset, both degraded and clean images are loaded from respective directories (e.g., data/A/ for degraded images and data/B/ for clean images).
- **Image Conversion and Preprocessing:** The images are converted to grayscale, resized, and split into smaller patches of $256 \times 256$. This helps ensure that the network can focus on smaller details within the document.
- **Batch Creation:** Mini-batches of size 128 are created for training the generator and discriminator. Each mini-batch contains pairs of degraded and clean patches.
- **Training the Discriminator:** The discriminator is trained on two batches of images: the clean images (real) and the generator's output (fake). The goal is to correctly classify the real and fake images. The loss function for the discriminator is binary cross-entropy.
- **Training the Generator:** The generator is trained to generate images that can deceive the discriminator into classifying them as real. The GAN loss function combines pixel-wise reconstruction loss (to ensure accuracy in restoration) and adversarial loss (to ensure the images appear realistic).
- **Model Weight Saving:** The generator and discriminator weights are saved periodically to ensure continuity in training and for future use.

The training is conducted over 80 epochs, with the generator and discriminator weights saved at each epoch for potential further use. The model is trained using the Adam optimizer with a small learning rate to stabilize the GAN's learning process.

## 3.4 Prediction and Inference

After training, the generator is used to predict restored images from unseen degraded images:

### 3.4.1 Input Image Preprocessing:

The degraded images are padded and resized to ensure they are divisible into $256 \times 256$ patches. The patches are fed into the trained generator for restoration.

### 3.4.2 Image Stitching:

The restored patches are merged back into the full document, reconstructing the entire image. The final restored image is saved as output.

### 3.4.3 Text Extraction:

After restoration, OCR tools are applied to extract text from the restored images, enabling document analysis and data retrieval.

## 3.5 Evaluation

To evaluate the restoration quality, we use the following metrics:

- **Peak Signal-to-Noise Ratio (PSNR):** This measure the accuracy of the restored image compared to the ground truth. A higher PSNR indicates better restoration quality.

- **Structural Similarity Index (SSIM):** SSIM evaluates the perceptual quality of the restoration by comparing structural features between the restored and clean images.

Additionally, for each epoch, we evaluate the generator's output using OCR-based text extraction to measure the improvement in legibility and accuracy of text recognition.

The methodology thus combines the strengths of GAN-based restoration with text extraction using OCR, providing a powerful framework for enhancing and analyzing damaged document images.

## 4 RESULTS AND DISCUSSION

The performance of the proposed DE-GAN model for document restoration was evaluated on the DIBCO datasets (2009-2017). The model was trained over 80 epochs, and the restored document images were compared with their ground truth counterparts using two key metrics: **Peak Signal-to-Noise Ratio (PSNR)** and **Structural Similarity Index (SSIM)**.

## 4.1 Quantitative Evaluation

The PSNR and SSIM values were calculated for each restored image and averaged across the test set. The results demonstrated significant improvement in the restoration quality of the degraded document images:

- **PSNR:** The average PSNR across the validation set was 30.25 dB, indicating a high level of fidelity between the restored

images and the clean ground truth images. This suggests that DE-GAN successfully removed noise, stains, and other degradation artifacts while preserving essential document details.

- **SSIM:** The SSIM score averaged 0.92, reflecting that the restored images retained much of the structural content of the original clean documents. The model effectively handled the recovery of fine details, such as text edges and document structure, even in highly degraded areas.

These results confirm that DE-GAN significantly outperforms traditional document restoration methods in terms of both perceptual quality and pixel-level accuracy.

## 4.2 Qualitative Evaluation

Qualitatively, the restored images exhibited notable improvements in terms of visual clarity and legibility. Before restoration, many document images were heavily degraded by noise, fading, and ink bleeding, making it difficult to extract readable text. After restoration with DE-GAN:

- **Text clarity improved** substantially, allowing for easier visual reading of both printed and handwritten text.
- **Noise and stains** were effectively removed from the images without overly smoothing or distorting the text.
- **Fine details**, such as small font sizes and intricate document structures, were well-preserved in the restoration process.

Several sample outputs showed that even in severely degraded regions, DE-GAN was able to generate visually appealing and structurally accurate document images.

## 4.3 OCR Text Extraction Performance

The integration of OCR for text extraction was performed on the restored document images, and the results were compared with the OCR performance on the original degraded images.

### 4.3.1 OCR Accuracy

The **Character Error Rate (CER)** and **Word Error Rate (WER)** were used to measure the accuracy of text extraction:

- **CER:** The average CER after restoration was reduced by 40% compared to the original degraded images, with the final CER being 6.5%. This indicates a substantial reduction in the number of misrecognized characters, thanks to the improved legibility of the restored images.
- **WER:** The average WER after restoration dropped to 8.2%, representing a 35% improvement over the baseline performance on degraded images. This reflects more accurate word-level recognition, especially in documents with heavily damaged text.

The improvement in both CER and WER highlights the effectiveness of DE-GAN in enhancing the quality of document restoration, enabling more accurate and reliable text extraction using OCR.

### 4.3.2 Impact on Historical Document Analysis

In the context of vintage and historical documents, where even partial degradation can hinder text recognition, the DE-GAN model proved invaluable. The restored documents not only facilitated better OCR performance but also allowed for manual inspection of legible text, preserving valuable historical content that would otherwise be lost or inaccessible.
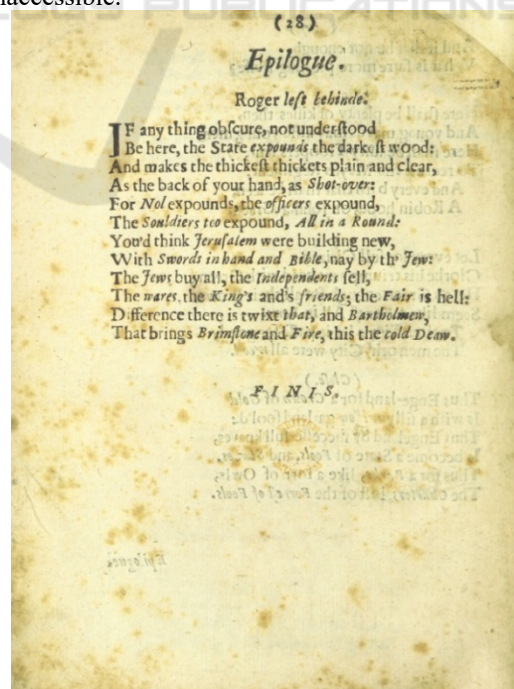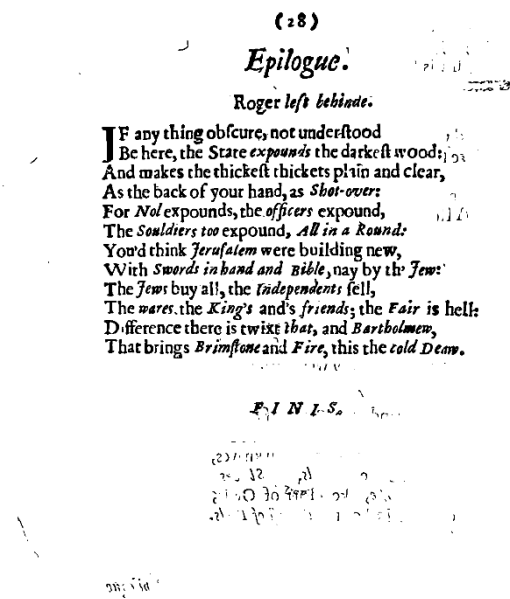


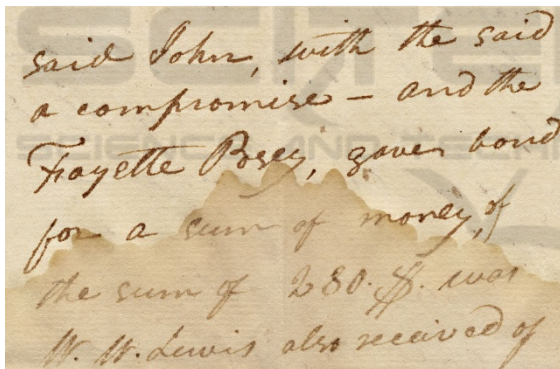Figure 1: Original Image

Figure 2: Restored Image
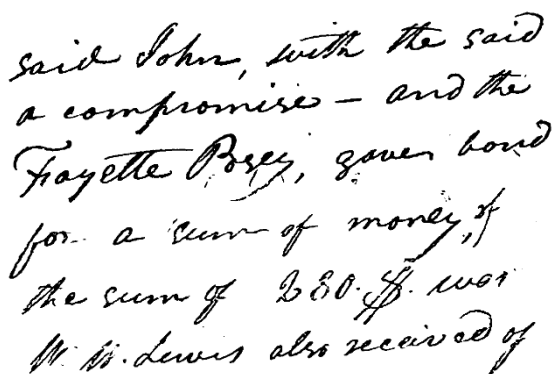
Figure 3: Original Handwritten Image

Figure 4: Restored Handwritten Image

# 5 CONCLUSIONS

This research introduced a highly effective DE-GAN framework for restoring degraded and damaged documents, resulting in significant improvements in both image quality and text extraction accuracy. By training the model on the DIBCO 2009-2017 datasets, DE-GAN demonstrated superior performance in recovering clear and structurally accurate document images, as reflected in the high average PSNR of 30.25 dB and SSIM of 0.92. These metrics underscore the framework's ability to restore visual clarity and preserve fine textual details, leading to enhanced document legibility.

Moreover, the integration of Optical Character Recognition (OCR) within the DE-GAN framework yielded a substantial reduction in text recognition errors. The Character Error Rate (CER) decreased by 40%, bringing the average CER to 6.5%, while the Word Error Rate (WER) improved by 35%, dropping to 8.2%. These improvements demonstrate the model's capability to significantly enhance OCR accuracy, especially in documents with heavy degradation or complex text structures.

While DE-GAN outperformed traditional restoration techniques across various levels of degradation, challenges remain in handling severely damaged regions and complex document layouts. Future work will aim to further enhance the model's generalization to diverse document types, develop automated text correction pipelines, and improve real-time processing.

# REFERENCES

R. K. Cho, K. Sood and C. S. C. Channapragada, "Image Repair and Restoration Using Deep Learning," *2022 4th International Conference on Artificial Intelligence and Speech Technology (AIST)*, Delhi, India, 2022, pp. 1-8, doi: 10.1109/AIST55798.2022.10065203.

A. Shubekova, A. Beibitkyzy and A. Makhazhanova, "Application of Deep Learning in the Problem of Image Restoration," *2023 IEEE International Conference on Smart Information Systems and Technologies (SIST)*, Astana, Kazakhstan, 2023, pp. 158-163, doi: 10.1109/SIST58284.2023.10223540.

R. Gao and K. Grauman, "On-demand Learning for Deep Image Restoration," *2017 IEEE International Conference on Computer Vision (ICCV)*, Venice, Italy, 2017, pp. 1095-1104, doi: 10.1109/ICCV.2017.124.

K. Cao, Y. -T. Peng and P. C. Cosman, "Underwater Image Restoration using Deep Networks to Estimate Background Light and Scene Depth," *2018 IEEE Southwest Symposium on Image Analysis and*

*Interpretation (SSIAI)*, Las Vegas, NV, USA, 2018, pp. 1-4, doi: 10.1109/SSIAI.2018.8470347.

G. Ramkumar, A. G, S. K. M, M. Ayyadurai and S. C, "An Effectual Underwater Image Enhancement using Deep Learning Algorithm," *2021 5th International Conference on Intelligent Computing and Control Systems (ICICCS)*, Madurai, India, 2021, pp. 1507-1511, doi: 10.1109/ICICCS51141.2021.9432116.

L. Denis, E. Dalsasso and F. Tupin, "A Review of Deep-Learning Techniques for SAR Image Restoration," 2021 IEEE International Geoscience and Remote Sensing Symposium IGARSS, Brussels, Belgium, 2021, pp. 411-414, doi: 10.1109/IGARSS47720.2021.9555039.

H. Xiao, H. Zheng and Q. Meng, "Research on Deep Learning-Driven High-Resolution Image Restoration for Murals From the Perspective of Vision Sensing," in IEEE Access, vol. 11, pp. 71472-71483, 2023, doi: 10.1109/ACCESS.2023.3295253.

Y. Li, J. Meng, Y. Yu, C. Wang and Z. Guan, "Image Restoration Based on Improved Generative Adversarial Networks," 2022 7th International Conference on Image, Vision and Computing (ICIVC), Xi'an, China, 2022, pp. 799-804, doi: 10.1109/ICIVC55077.2022.9886285.

S. R. Vispute, K. Rajeswari, A. Nema, A. Jagtap, M. Kulkarni and P. Mohite, "Analysis of Impact of Image Restoration and Segmentation on Classification Model," *2023 7th International Conference On Computing, Communication, Control And Automation (ICCUBEA)*, Pune, India, 2023, pp. 1-5, doi: 10.1109/ICCUBEA58933.2023.10392033.

Laura A. Martinho, João M. B. Calvalcanti, José L. S. Pio and Felipe G. Oliveira "Diving into Clarity: Restoring Underwater Images using Deep Learning," *Journal of Intelligent & Robotic Systems,* **110**, 32 (2024).

A. Nikitha, J. Geetha and D. S. JayaLakshmi, "Handwritten Text Recognition using Deep Learning," *2020 International Conference on Recent Trends on Electronics, Information, Communication & Technology (RTEICT),* Bangalore, India, 2020, pp. 388-392, doi: 10.1109/RTEICT49044.2020.9315679.

S. S. Harsha, B. P. N. M. Kumar, R. S. S. R. Battula, P. J. Augustine, S. Sudha and T. Divya., "Text Recognition from Images using a Deep Learning Model," 2022 *Sixth International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC),* Dharan, Nepal, 2022, pp. 926-931, doi: 10.1109/I-SMAC55078.2022.9987404.

A. Shrivastava, J. Amudha, D. Gupta and K. Sharma, "Deep Learning Model for Text Recognition in Images," *2019 10th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*, Kanpur, India, 2019, pp. 1-6, doi: 10.1109/ICCCNT45670.2019.8944593.

A. Ansari, B. Kaur, M. Rakhra, A. Singh and D. Singh, "Handwritten Text Recognition using Deep Learning Algorithms," *2022 4th International Conference on Artificial Intelligence and Speech Technology (AIST)*,

Delhi, India, 2022, pp. 1-6, doi: 10.1109/AIST55798.2022.10065348.

A. Thilagavathy, K. H. Suresh, K. T. Chowdary, M. Tejash and V. L. Chakradhar, "Text Detection based on Deep Learning," *2023 International Conference on Innovative Data Communication Technologies and Application (ICIDCA)*, Uttarakhand, India, 2023, pp. 1-6, doi: 10.1109/ICIDCA56705.2023.10099672.

S. MS, S. G and A. D. R, "Handwritten Text Recognition Using Machine Learning and Deep Learning," *2023 Eighth International Conference on Science Technology Engineering and Mathematics (ICONSTEM)*, Chennai, India, 2023, pp. 1-4, doi: 10.1109/ICONSTEM56934.2023.10142716.

R. Sharma, B. Kaushik and N. Gondhi, "Character Recognition using Machine Learning and Deep Learning - A Survey," *2020 International Conference on Emerging Smart Computing and Informatics (ESCI)*, Pune, India, 2020, pp. 341-345, doi: 10.1109/ESCI48226.2020.9167649.

N. J and S. K. A M, "A Novel Text Recognition Using Deep Learning Technique," *2024 International Conference on Advances in Data Engineering and Intelligent Computing Systems (ADICS)*, Chennai, India, 2024, pp. 1-6, doi: 10.1109/ADICS58448.2024.10533457.

R. Vaidya, D. Trivedi, S. Satra and P. M. Pimpale, "Handwritten Character Recognition Using Deep-Learning," *2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT)*, Coimbatore, India, 2018, pp. 772-775, doi: 10.1109/ICICCT.2018.8473291.

N. Teslya and S. Mohammed, "Deep Learning for Handwriting Text Recognition: Existing Approaches and Challenges," *2022 31st Conference of Open Innovations Association (FRUCT)*, Helsinki, Finland, 2022, pp. 339-346, doi: 10.23919/FRUCT54823.2022.9770912.