

Super Resolution Image Denoiser: Toward Real-world Noise Removal

P. Mohit Harsh

*Electronics and Communication Engineering
Institute of Aeronautical Engineering
Hyderabad, India
21951a04a3@iare.ac.in*

V. Teja Vardhan

*Electronics and Communication Engineering
Institute of Aeronautical Engineering
Hyderabad, India
21951a04q8@iare.ac.in*

V. Vishnu Vardhan

*Electronics and Communication Engineering
Institute of Aeronautical Engineering
Hyderabad, India
21951a04q9@iare.ac.in*

J. Manoj Naidu

*Electronics and Communication Engineering
Institute of Aeronautical Engineering
Hyderabad, India
21951a04r4@iare.ac.in*

Dr. V. Padmanabha Reddy

*Electronics and Communication Engineering
Institute of Aeronautical Engineering
Hyderabad, India
v.padmanabhareddy@iare.ac.in*

Abstract—In recent years, deep learning has significantly advanced image denoising techniques, with methods such as transformers and GANs pushing the boundaries of performance. However, these state-of-the-art models are typically resource-intensive, requiring vast computational power, large datasets, and extensive training time, which also hampers real-time inference speed. To address these challenges, we propose the Super-Resolution Image Denoiser Network (SRIDNet), a lightweight yet highly effective model designed to achieve near state-of-the-art results with a fraction of the data and computational resources. SRIDNet significantly reduces the model size while maintaining competitive performance in image denoising and super-resolution tasks. We evaluated our model on multiple datasets, where it demonstrated high PSNR scores and the fastest inference times compared to existing models. Our results confirm that SRIDNet offers an efficient solution for real-time image denoising, making it a practical alternative to more cumbersome models without compromising on quality.

Index Terms—Convolution Neural Network (CNN), Adversarial Neural Network (GAN), SRIDNet, PSNR, Transformers.

I. INTRODUCTION

Image denoising is a critical preprocessing step in many computer vision and image processing applications, aimed at restoring corrupted images by reducing noise while preserving important structural details. Traditional denoising techniques, such as Gaussian filtering and wavelet transforms, have long been employed for this task but are often limited by their inability to adapt to complex and diverse noise patterns encountered in real-world images. In recent years, the rise of deep learning has revolutionized the field of image denoising, offering powerful data-driven solutions capable of learning

intricate noise distributions and producing visually appealing results.

Convolutional Neural Networks (CNNs), in particular, have become the dominant paradigm for image denoising, leveraging their ability to capture local spatial hierarchies through multiple layers of convolutional filters. These models have significantly outperformed classical approaches, achieving state-of-the-art results by learning from large datasets of noisy and clean image pairs. Recent advancements have also introduced more sophisticated architectures such as Generative Adversarial Networks (GANs) and Transformer models, further improving denoising performance by capturing global dependencies and generating more realistic outputs.

Despite their impressive results in controlled environments, these state-of-the-art models often come with high computational costs, extensive memory requirements, and the need for large training datasets. These factors limit their practicality in real-world applications, particularly in resource-constrained environments where hardware and data availability are restricted.

To address these challenges, we propose the Super Resolution Image Denoiser Network (SRIDNet), a novel approach that strikes a balance between denoising performance and computational efficiency. SRIDNet is designed to deliver competitive results with significantly lower hardware demands and reduced model complexity, making it suitable for deployment in practical settings. By integrating super-resolution techniques with an efficient denoising mechanism, SRIDNet provides a robust solution for noise reduction without sacrificing speed or resource efficiency. This paper presents the architecture,

training methodology, and performance evaluation of SRID-Net, showcasing its advantages in terms of both accuracy and efficiency across various datasets and noise levels.

II. LITERATURE SURVEY

Deep convolutional neural networks (CNNs) have made substantial progress in the last ten years in high-level vision tasks such as object segmentation, motion analysis, and visual recognition. CNNs have been used more recently for low-level vision tasks like compression artifact reduction, image denoising, and super-resolution (SR). In these tasks, CNNs are trained to map low-quality images to high-quality outputs, usually with the goal of minimizing artifacts or removing noise. While "deeper is better" is a widely accepted principle in high-level vision tasks—evidenced by networks like VGG, GoogleNet, and ResNet achieving substantial breakthroughs—this principle has not shown as much impact in low-level vision tasks. Despite the use of networks with 20 to 30 layers, such as DnCNN [1] and RED-Net, the performance gains in low-level tasks have been modest compared to earlier methods. This is because low-level vision tasks rely more on pixel-level features, where depth is less crucial. Instead, statistical priors, like non-local similarity or pixel distribution patterns (e.g., Gaussian noise), play a key role in enhancing the accuracy of these tasks, offering a more effective solution to image degradation issues.

A. Existing Work

In CBDNet [2], an asymmetric loss function was employed to enhance the model's ability to generalize to real-world noise scenarios, while also facilitating convenient interactive denoising. The model demonstrated PSNR scores of 30.78 on the SIDD dataset [5] and 38 on the DnD dataset. Processing a 512x512 image takes approximately 0.4 seconds.

In their work on RIDNet [3], the authors introduce a CNN-based denoising model specifically designed for both synthetic noise and real-world noisy images. The model's architecture consists of four Enhanced Attention Mechanism (EAM) blocks, where most convolutional layers have a kernel size of 3x3, except for the final layer in the enhanced residual block and feature attention units, which utilize a 1x1 kernel. The model processes a 512x512 image in approximately 0.2 seconds during evaluation and achieved PSNR scores of 31.38, 39.23, and 38.71 on the BSD68 [4], DnD [6], and SIDD [5] datasets, respectively.

The NBNet [7] architecture is built upon a modified UNet framework. NBNet incorporates four encoder and decoder stages, where feature maps are downsampled using strided convolutions in the encoder and upsampled using deconvolutions in the decoder. It achieved PSNR scores of 29.16, 39.62 and 39.75 on BSD68 [4], DnD [6] and SIDD [5] datasets respectively.

DANet [2] is the latest architecture introduced for real-world image denoising tasks. Its core concept revolves around unfolding in-camera processing pipelines or learning the noise distribution through a generative adversarial network (GAN).

This architecture achieved PSNR scores of 39.25 and 39.79 on the SIDD and DnD benchmark datasets, respectively.

B. Evaluation Metrics

In image processing, Structural Similarity Index Measure (SSIM), Mean Squared Error (MSE), and Peak Signal-to-Noise Ratio (PSNR) are frequently used metrics to evaluate the quality of images, especially when determining how similar a processed image is to its reference image.

1) *Mean Squared Error (MSE)*: MSE is a pixel-based metric that calculates the average of the squared differences between corresponding pixels of two images—typically, an original (reference) image and a processed (distorted or denoised) image.

$$MSE = \frac{1}{m \times n} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} (I(i, j) - K(i, j))^2$$

Where:

- The value of the pixel at location (i, j) in the original image is represented by $I(i, j)$.
- The pixel value at location (i, j) in the processed image is denoted by $K(i, j)$.
- The image dimensions are m and n .

However, MSE doesn't take human visual perception into account and is sensitive to large differences, even if small changes are hard to perceive.

2) *Peak Signal-to-Noise Ratio (PSNR)*: PSNR is a more perceptually aligned metric that expresses the ratio between the maximum possible pixel value of the image and the error (noise) introduced by the processing, measured via MSE.

$$PSNR = 10 \times \log_{10} \left(\frac{MAX_I^2}{MSE} \right)$$

Where:

- MAX_I is the maximum possible pixel value (255 for an 8-bit image).
- MSE is the mean squared error between the two images.

Interpretation:

- Higher PSNR indicates that the processed image is of better quality and closer to the original.
- A PSNR above 30 dB is generally considered acceptable for most image processing tasks, though in applications like medical imaging, higher values are preferred.

PSNR is easy to compute and gives a high-level indication of image quality, but it still does not correlate well with human visual perception.

C. Drawbacks of existing methods

1) *High Computational Complexity and Resource Demand*: State-of-the-art denoising models typically consist of millions of trainable parameters, leading to significant computational demands. While these models perform exceptionally well in controlled laboratory environments, their need for vast datasets and computational resources makes them impractical for many real-world applications.

- *Concatenation with Input*: A skip connection is established where the original input image is concatenated with the downsampled version of the upscaled image, producing an intermediate feature map of size $(112 \times 112 \times 131)$. This fusion ensures the model retains essential features from the initial input.
- *Downsampling and Upsampling*:
 - The image is further downsampled using another convolutional layer with 64 filters, followed by another *MaxPooling2D* operation that reduces the spatial size to (56×56) .
 - An *UpSampling2D* layer increases the resolution back to (112×112) .
- *Concatenation with Intermediate Features*: Skip connections are introduced again by concatenating the upsampled feature map with previous intermediate layers, producing a feature map of size $(112 \times 112 \times 195)$.
- *Final Convolutions*: The concatenated feature map is passed through a convolutional layer with 128 filters and finally reduced to 3 channels (RGB) using a *Conv2D* layer with 3 filters. The output is the denoised image, reconstructed to its original resolution of (112×112) .

D. Working Principle

- The architecture leverages the *Super-Resolution Block* to enhance the resolution of the input image from (112×112) to (224×224) . This block employs a combination of *convolutional layers*, *residual connections*, and *depth-to-space transformation* to ensure that the upscaled image retains fine details without introducing significant artifacts.
- Following this, the *Denoiser Block* takes over, utilizing downsampling and upsampling techniques, combined with skip connections, to effectively remove noise from the upscaled image. The final result is a clean, denoised image at the original resolution of (112×112) .
- *Skip connections* play a crucial role in both blocks by preserving critical features from the earlier stages and preventing information loss during downsampling and upsampling.

E. Summary of Parameters

- *Total Parameters*: 454,902
- *Trainable Parameters*: 454,774
- *Non-trainable Parameters*: 128

This architecture is highly optimized for the task of image denoising, leveraging the efficiency of the *Super-Resolution Block* to work with smaller input patches, significantly reducing the amount of data and time required for training and inference. Furthermore, the *Denoiser Block* ensures that the model can effectively remove noise while retaining critical image features, providing high-quality denoised images as output.

IV. IMPLEMENTATION

A. Data Preprocessing

1) *SIDD* [5]: In our study, we utilized images from the SIDD [5] dataset, resizing them to dimensions of $(2576, 1456, 3)$. From each image, we extracted patches with a resolution of $(112, 112, 3)$. For training, we randomly selected 7,000 patches, while 3,000 patches were allocated for validation. The final 10 images from the dataset were reserved for testing, which, after patching, yielded 2,990 test patches.

2) *Urban100* [9]: In our study, we utilized the Urban100 [9] dataset to extract a total of 80 images for the training set. These images were cropped to a size of $560 \times 560 \times 3$ pixels. Subsequently, we generated patches of size $112 \times 112 \times 3$ from each image, resulting in a total of 2,000 patches designated for training with a validation split of 0.1. For the testing phase, we selected 20 remaining images from the dataset, which were similarly cropped to $560 \times 560 \times 3$ pixels. Following the patch generation process, we obtained 500 patches of size $112 \times 112 \times 3$ for testing purposes.

B. Model Training

The Super-Resolution Image Denoiser Network (SRIDNet) was trained using the TensorFlow and Keras frameworks on image patches of size $(112, 112, 3)$. The training process was carried out in three stages, each with progressively reduced learning rates to enhance model performance and stability. The model was first trained using a learning rate of 0.001 across 20 epochs. After then, there were 10 further epochs with a further decreased learning rate of 0.00001, and 20 more epochs at a lowered learning rate of 0.0001. Peak Signal-to-Noise Ratio (PSNR) was utilized as the assessment metric, and the Adam optimizer was utilized to minimize the mean squared error (MSE) loss function. The batch size was set to 16 for all training stages. On the Urban100 dataset, the model required approximately 45 seconds per epoch, while on the SIDD dataset, it required 130 seconds per epoch. These timings reflect the computational complexity and dataset size differences between the two datasets.

V. RESULTS AND COMPARISONS

The proposed Super-Resolution Image Denoiser Network (SRIDNet) was evaluated on the Urban100 [9] dataset, with the last 10 images used as the test set. The evaluation process was conducted with TensorFlow, and inference was completed in an average time of 20-24 ms per image of size $560 \times 560 \times 3$ at noise level ($\sigma = 15$). The model's performance was assessed using Peak Signal-to-Noise Ratio (PSNR) as the evaluation metric. The results indicate that SRIDNet achieved an average PSNR of 34.23 dB across the test images, with a maximum PSNR of 35.82 dB.

The proposed Super-Resolution Image Denoiser Network (SRIDNet) was evaluated on the Urban100 dataset and compared against state-of-the-art denoising models, including DnCNN [1], FFDNet [12], DRUNet [13], and SwinIR [14], across different noise levels ($\sigma = 15, 25, 50$). The evaluation

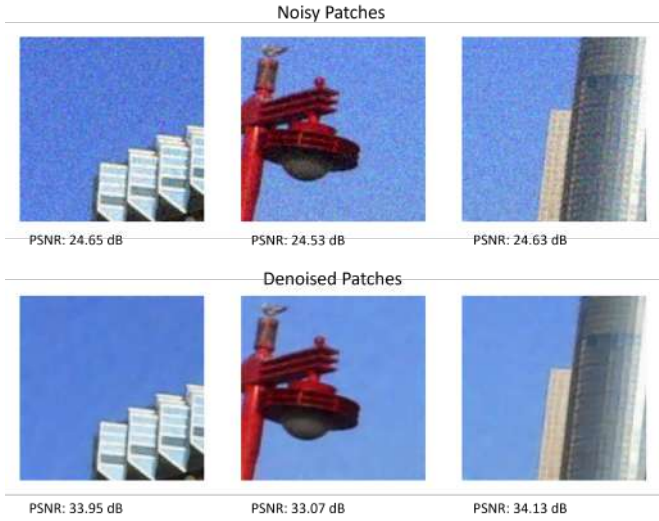


Fig. 4. PSNR results on $112 \times 112 \times 3$ image patches ($\sigma = 15$) from Urban100 dataset

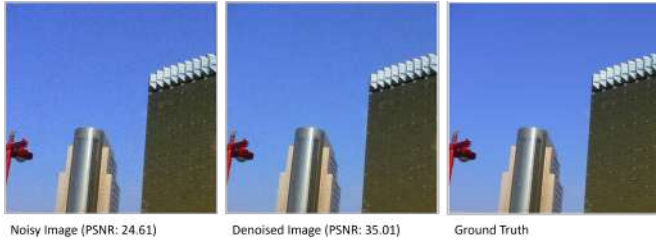


Fig. 5. PSNR results on an entire image ($\sigma = 15$) of size $560 \times 560 \times 3$ from Urban100 dataset

metric is Peak Signal-to-Noise Ratio (PSNR) as presented in table 1.

The proposed model was evaluated on the SIDD dataset and benchmarked against several state-of-the-art denoising models, including DANet [16], VDN [15], RIDNet [3], and CBDNet [2]. The model achieved an average PSNR of 39.5 on the SIDD dataset, with an average inference time of 95 ms for images sized 2576×1456 , which were divided into 299 patches of size 112×112 .

Figures 6 and 7 illustrate the denoising results of low-light and bright-light images from the SIDD [5] dataset, respectively, using $(112 \times 112 \times 3)$ image patches.

Table II presents a comparison of denoising performance across different models on the SIDD [5] dataset, showing that our model outperformed others by achieving a higher average PSNR. While DnCNN [1] and FFDNet [12]—models of similar size—performed well on images with additive Gaussian white noise (AGWN) from Urban100 [9] dataset, they underperformed on real-world photographic images from the SIDD [5] dataset. In contrast, our model demonstrated robust performance across both datasets with greater efficiency.

Table III compares the inference times of the models evaluated in this study. While our model is faster than all models except FFDNet [12], it offers more consistent results

Model	σ	PSNR(dB)
DnCNN	15	32.98
	25	30.81
	50	27.59
FFDNet	15	33.83
	25	31.40
	50	28.05
DRUNet	15	34.81
	25	32.60
	50	29.61
SwinIR	15	35.13
	25	32.90
	50	29.82
ours	15	34.23
	25	32.01
	50	28.67

TABLE I
URBAN100 IMAGE DENOISING RESULTS AND COMPARISON WITH EXISTING MODELS

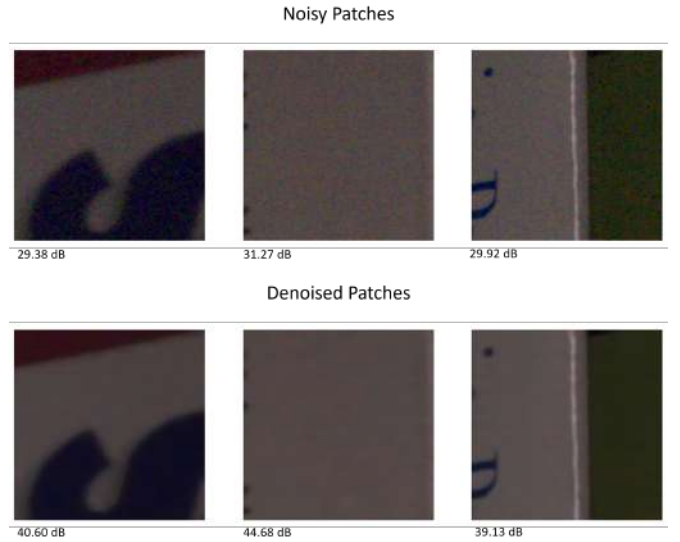


Fig. 6. Denoising results on SIDD dataset images

on both the Urban100 and SIDD [5] datasets, where FFDNet [12] struggles with real-world images.

Model	PSNR (dB)
CBDNet	30.78
RIDNet	38.71
VDN	39.28
DANet	39.47
DnCNN	26.21
FFDNet	29.20
Ours	39.50

TABLE II
COMPARISON OF RESULTS WITH STATE OF THE ART MODELS ON SIDD [5] DATASET

VI. CONCLUSION

In conclusion, this research presents SRIDNet, a novel deep learning architecture for image denoising that effectively addresses the limitations of current state-of-the-art mod-



Fig. 7. Denoising results on SIDD dataset images

Model	Inference Time (s)
DnCNN	0.0314
FFDNet	0.0071
DRUNet	0.0733
CBDNet	0.4
RIDNet	0.2
Ours	0.0208

TABLE III

COMPARISON OF INFERENCE TIME WITH DIFFERENT MODELS.

els, such as high computational demands, extensive training data requirements, and slow inference times. By designing a lightweight model that integrates super-resolution and denoising tasks, SRIDNet achieves competitive results with significantly reduced resource consumption. Extensive testing on both synthetic (AGWN) and real-world noisy images from the Urban100 and SIDD datasets demonstrates the model's robustness and efficiency. With PSNR scores of 34.23 on Urban100 and 39.50 on SIDD, SRIDNet delivers near state-of-the-art performance while maintaining the smallest model size and fastest inference time. This balance of efficiency and effectiveness positions SRIDNet as a promising solution for practical image denoising applications, particularly in environments with limited computational resources.

VII. ACKNOWLEDGMENT

We would like to extend our heartfelt gratitude to those who made this project possible through their constant guidance and support. Our deepest appreciation goes to our guide, Dr. V. Padmanabha Reddy, Assistant Professor, for his invaluable guidance and continuous cooperation throughout the project. We also want to thank the teaching and non-teaching staff for their support, and express our deepest gratitude to our parents, friends, and well-wishers for their assistance in completing this project report successfully.

REFERENCES

- [1] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE Transactions on Image Processing*, 26:3142–3155, 2017.
- [2] S. Guo, Z. Yan, K. Zhang, W. Zuo, L. Zhang: Toward Convolutional Blind Denoising of Real Photographs. In: *CVPR* (April 2019).
- [3] S. Anwar, N. Barnes.: Real Image Denoising with Feature Attention. In: *ICCV* (March 2020).
- [4] Stefan Roth and Michael J Black. Fields of experts. *IJCV*, 2009. 1, 2, 5, 6, 7
- [5] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *CVPR*, 2018. 5, 8
- [6] Tobias Plotz and Stefan Roth. Benchmarking denoising algorithms with real photographs. *arXiv preprint arXiv:1707.01313*, 2017. 5, 7
- [7] S. Chengl, Y. Wangl, H. Huang, D. Liu, H. Fan and S. Liu.: NBNNet: Noise Basis Learning for Image Denoising with Subspace Projection. In: *CVPR* (May 2021).
- [8] Sara, U. , Akter, M. and Uddin, M. (2019) Image Quality Assessment through FSIM, SSIM, MSE and PSNR—A Comparative Study. *Journal of Computer and Communications*, 7, 8-18. doi: 10.4236/jcc.2019.73002.
- [9] J.-B. Huang, A. Singh, and N. Ahuja, "Single image super-resolution from transformed self-exemplars," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, June 2015.
- [10] V. Bychkovsky, S. Paris, E. Chan, and F. Durand, "Learning photographic global tonal adjustment with a database of input/output image pairs," in *Proc. 24th IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2011.
- [11] K. Ma, Z. Duanmu, Q. Wu, Z. Wang, H. Yong, H. Li, and L. Zhang, "Waterloo Exploration Database: New challenges for image quality assessment models," *IEEE Transactions on Image Processing*, vol. 26, no. 2, pp. 1004-1016, Feb. 2017.
- [12] Kai Zhang, Wangmeng Zuo, and Lei Zhang. Ffdnet: Toward a fast and flexible solution for cnn-based image denoising. *TIP*, 2018. 1, 2, 5, 6, 7, 8
- [13] Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 1, 2, 7, 8
- [14] J. Liang, J. Cao, G. Sun, K. Zhang, L. Van Gool, and R. Timofte, "SwinIR: Image Restoration Using Swin Transformer," *arXiv preprint arXiv:2108.10257*, 2021.
- [15] Z. Yue, H. Yong, Q. Zhao, D. Meng, and L. Zhang, "Variational denoising network: Toward blind noise modeling and removal," in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 1690-1701. Available: <http://papers.nips.cc/paper/8446-variational-denoising-network-toward-blind-noise-modeling-and-removal.pdf>.
- [16] Z. Yue, Q. Zhao, L. Zhang, and D. Meng, "Dual adversarial network: Toward real-world noise removal and noise generation," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, Aug. 2020.