



Image Super Resolution Enhancement Using A Hybrid Model Framework

Uday Sharma^{1*}, Kanika Mahajan¹, Sanjana Kharbanda¹, Shubham Kathuria¹ and Swapnil Kaushal¹

¹Dept. of Computer Science and Engineering, JIMS Engineering Management Technical Campus, Greater Noida, India

*uday69510@gmail.com

kani.mahaji@gmail.com, sanjanak0607@gmail.com, shubhamkathuria06@gmail.com, swapnilkaushal.gn@jagannath.org

Abstract. The goal of Image Super-Resolution is to obtain high-resolution images from low-resolution images, but it is still difficult to obtain both structure accuracy and perceptual sharpness. CNN models have been shown to restore edges and textures, but results can often be overly smooth. Transformer architectures are able to model long-range spatial dependency effectively, but can lead to structural inconsistencies. Likewise, GAN-based image enhancement gives a high degree of perceptual realism, but can often produce undesirable artifacts when used alone. In response to these limitations, this paper proposes a new approach to SR by presenting a hybrid Super-Resolution framework based on a training pipeline that includes Convolutional Neural Networks, Transformer encoders, and a Generative Adversarial Network. The process of training the model progresses through four stages: baseline CNN reconstruction, global context refinements from the Transformer-oriented architectural attention, adversarial learning for perceptual textures, and finally, a Fusion Network that adaptively blends outputs from the three models. The DIV2K dataset is employed for training and evaluation. The hybrid framework composites complementary variations of the individual architectures in an effort to achieve improvements in detail preservation, structural coherence, and perceptual realism compared to their respective baselines employing CNN-only, Transformer-only, or GAN-only approaches to super-resolution. Overall, we believe the proposed framework has great potential for real-world uses, such as medical imaging, satellite observation, and visual restoration applications, when accuracy and reliability are paramount.

Keywords: Super Resolution, Convolutional Neural Network (CNN), Transformer, Generative Adversarial Network (GAN), Fusion Network, Hybrid Framework.

1 Introduction

High resolution images are important in areas such as medical diagnosis, satellite monitoring, surveillance, and digital restoration of media. Nonetheless, imaging sensors, environmental conditions, or storage/transmission limitations typically result

in low resolution images that lack fine texture details, sharp edges, and structural coherence. Image Super-Resolution is the process of reconstructing a high-resolution image from a low-resolution input, while achieving perceptual realism and structural integrity has remained difficult.

Conventional interpolation-based techniques increase the size of an image but do not recover lost details. Convolutional Neural Networks have improved image quality partially by learning spatial feature representations [1], [5], [6]. However, due to their limited receptive fields, CNNs may output overly smooth images [2], [4]. Transformer-based approaches facilitate learning the global spatial relationships of pixels applied through self-attention method, but they will struggle with fine textures as they are used in isolation [12], [13], [14], [28], [32]. Generative Adversarial Networks (GANs) enhance visual sharpness and clarity but may create artefacts or unnatural textures if the contrastive loss is not appropriately calibrated [2], [3], [11], [30], [31].

This indicates that there is no one model that can reliably build high-quality images. As such, we believe it is beneficial to build a hybrid that pulls together the strengths of multiple models; CNN, Transformer, and GAN. In this work, we propose a hybrid Super-Resolution framework that contains a CNN, Transformer, and GAN model that are trained sequentially, allowing each to specialize prior to a Fusion Network sitting on top of the all three models to combine their outputs. We aim to provide improved clarity, structure consistency, and realism in texture restoration comparisons to three baseline models. We utilize the DIV2K dataset to train and evaluate this work [20], and we will compare and place our hybrid technique against current state-of-the-art super-resolution improvements [28], [29], [30], [31], [32], [33].

2 Literature Review

The founders of early Image Super-Resolution relied on interpolation schemes which upscaled an image, but didn't recover high-frequency signals, leading to blurry outputs, such as those that might result from bicubic upscaling [1]. Later, Super-Resolution based on deep networks introduced the use of Convolutional Neural Networks (CNNs); widely used models include VDSR, EDSR, RDN, and RCAN [1], [4], [5], [6], all which have contributed to advancements in reconstruction methods by learning to map the local features from a low to high-resolution image. While CNNs are effective with respect to edge and spatial structure restoration, their limited receptive fields hinder their ability to consider long-range dependencies, which typically leads to smooth textures and representation of complex images [2].

Transformer based approaches solved this issue associated with limited receptive fields through the use of self-attention mechanisms to model global pixel dependencies. For example, methods such as Image Transformer, IPT, SwinIR, HAT and DAT have shown promising results for performance with respect to global structure and contextual consistency [6], [12], [13], [14], [15], [28], [32]. The drawback of these methods is that self-attention is often computationally expensive and while self-attention-based networks might perform better than CNNs, they often fail to reconstruct finer perceptual textures when trained using strictly pixel loss [10].

Generative Adversarial Networks also improved SR through adversarial and perceptual learning. Frameworks such as SRGAN, ESRGAN, ESRGAN+, Real-ESRGAN, and A-ESRGAN can produce sharper, more realistic details by pushing for high-frequency texture recovery [2], [3], [9], [10], [11], [30], [31]. However, when improperly trained, GANs can reproduce hallucinated context or introduce artefacts [17].

Recent works tested in various domains show that CNNs, Transformers, and GAN architectures each play a different role in SR; CNNs capture a local structure, Transformers capture global structure, and GANs improve perceptual realism [18], [19], [27], [29], [33]. This behavior of complementarity drove the movement to hybrid approaches for SR, whereby combining models can balance quality of reconstruction. Further work on texture transfer, attention-based fusion, and fine-tuning pre-trained models for SR help to verify the strengths of using multiple SR models [16], [21], [23], [28], [29], [30], [31], [32], [33]. Table 1 summarizes recent studies on AI-based image super-resolution techniques, including their objectives, outcomes, and limitations.

Table 1. Summary of Recent Studies on AI-Based Image Resolution Enhancement and Reconstruction Techniques.

Author	Title	Objective	Outcome	Limitation
Lim et al., 2017 [1]	EDSR: Enhanced Deep Residual Networks for SISR	Deeper residual CNN for improved SR accuracy.	Achieved state-of-the-art PSNR/SSIM on DIV2K.	Limited receptive field; lacks global context modeling.
Wang et al., 2018 [2]	ESRGAN: Enhanced Super-Resolution GAN	Enhance SRGAN using RRDB blocks and perceptual loss.	Achieved superior perceptual quality with stable GAN training.	May over-emphasize textures and distort structure.
Zhang et al., 2018 [4]	RCAN: Residual Channel Attention Network	Channel attention for improved feature selectivity in CNN SR.	Improved fine-detail reconstruction and boosted fidelity metrics.	Still constrained by local CNN receptive fields.
Liang et al., 2021 [12]	SwinIR: Image Restoration Using Swin Transformer	Window-based Transformers for efficient global context modeling.	Achieved strong global consistency and high SR accuracy.	Computationally expensive; may miss high-frequency textures.
Baghel et al., 2023 [16]	SRTransGAN	Combine Transformers and GANs for global context and perceptual quality.	Demonstrated enhanced fidelity and perceptual quality over standalone models.	Training complexity increases; risk of GAN instability.

3 Methodology

The proposed hybrid image super-resolution framework intends to combine Convolutional Neural Networks, Transformer attention mechanisms, and a Generative Adversarial Network, as well as to add a Fusion Network to combine the strengths of all three types of model for a complete image super-resolution architecture [1], [3], [12], [28], [30], [31], [32]. The system is trained in four phases in order to promote stability of optimization and balanced visual reconstruction [9], [16], [29], [33]. The different phases of the proposed model are outlines in Table 2. The overall workflow of the proposed hybrid super resolution framework is illustrated in Fig. 1.

Table 2. Breakdown of the phases of the model.

Phase	Section Title	What It Covers
Phase 1	Dataset Preparation	DIV2K, down sampling, normalization, augmentation, patch extraction
Phase 2	Baseline CNN Model	Local feature extraction + Pixel Shuffle reconstruction
Phase 3	Transformer Model	Global context modeling using attention
Phase 4	GAN Model	Perceptual sharpness + adversarial refinement
Phase 5	Fusion Network	Final combination of outputs
Phase 6	Training Strategy	Stage 1–4 training schedule

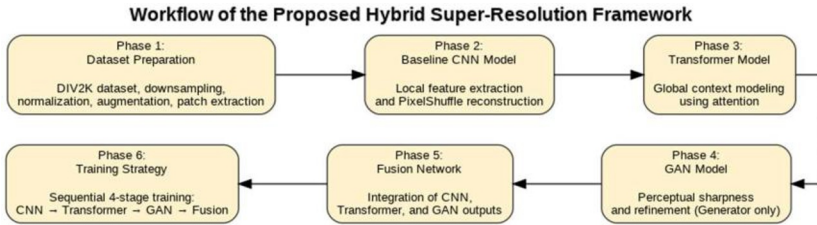


Fig. 1. Workflow Diagram of the model.

3.1 Dataset Preparation

The DIV2K dataset was selected as the primary benchmark due to its substantial diversity and for high-resolution ($\sim 2000 \times 2000$ pixels) images [20], [28]. Each high-resolution image was then down sampled via a bicubic interpolation operation at a $\times 4$ scale for low-resolution inputs:

$$I_{LR} = D_{\text{bicubic}}(I_{HR}, s = 4) \quad (1)$$

The images were all normalized to the range of [0,1] and augmented using horizontal/vertical flips and 90° rotations to facilitate generalization. A patch extraction applied for efficient training, with each 48×48 LR patch corresponding to its associated 192×192 HR patch. The dataset was then divided into 80% training, and a

test/validation split of 10% each to aid balanced learning of both fine texture and global structures. The generation process low resolution and high resolution image pairs is shown in Fig. 2.

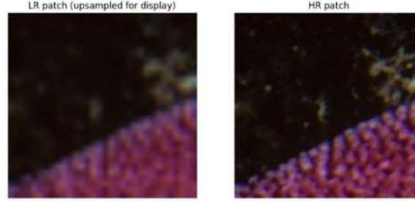


Fig. 2. Generation of LR–HR training pairs from the DIV2K dataset using bicubic down sampling at $\times 4$ scale.

Baseline CNN Model. - The initial step of the proposed framework introduces a Convolutional Neural Network that learns the underlying spatial mapping from each low-resolution input I_{LR} to its high-resolution counterpart I_{HR} , yielding a baseline super-resolved image [1], [5], [6], [30]. As the CNN focuses primarily on local feature extraction (e.g., edges, fine textures, and structural boundaries), it consists of several stacked 3×3 convolution layers activated with the ReLU function and subsequent residual blocks, following established SR designs [4]. The inherent property of residual learning retains spatial information and reduces gradient vanishing effects during training [6]. Each residual block refines features as follows:

$$F_{res}(x) = x + Conv_{3 \times 3}(\sigma(Conv_{3 \times 3}(x))) \quad (2)$$

where σ is the activation function (ReLU). The network employs Pixel Shuffle sampling to recover the high-resolution spatial grid; this process rearranges channel-wise feature maps to the spatial dimension without introducing checkerboard artifacts [2], resulting in finer detail reconstruction. Let F_{CNN} represent the CNN transformation, the baseline super-resolved output is:

$$I_{SR}^{CNN} = F_{CNN}(I_{LR}) \quad (3)$$

We train the CNN using L1 content loss in order to impose pixel-wise fidelity,

$$\mathcal{L}_{CNN} = \|I_{HR} - I_{SR}^{CNN}\|_1 \quad (4)$$

This stage provides a structurally consistent, smooth high-resolution image that serves as the foundation for later refinement by the Transformer and GAN components.

Transformer Based Global Context Modeling. - Although the CNN baseline is capable of restoring the local structures, its spatial receptive field is limited [4]. The second model leverages a patch-based Transformer architecture to model long-range dependencies across the image to capture global contextual relationships [12], [13], [14], [15], [28], [32]. The input image I_{LR} can be seen as a grid of non-overlapping spatial regions. Then the image is partitioned into 4×4 patches before being embedded

into a sequence of token embeddings which are generated via a convolutional patch embedding layer:

$$X_{\text{patch}}, (h, w) = \text{PatchEmbed}(I_{\text{LR}}) \quad (5)$$

It creates a token sequence of length hw of each patch representing a spatial region of the image. The token sequence is then processed through a stack of Transformer Encoder layers, applying Multi-Head Self-Attention and a feed-forward network for each layer:

$$\text{MHSA}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (6)$$

where Q, K, V are learned linear projections of the token embeddings. This allows each patch to attend to all other patches and gain global structured awareness [12], [13], [32]. Note that the input data is a sequential representation of the input image supporting the modeling of long-range features. Subsequent to processing through the Transformer Layers, the feature tokens are reshaped back into the spatial grid to once again represent the image and are passed through a ConvTranspose2D + Pixel Shuffle reconstruction module:

$$I_{\text{SR}}^{\text{TR}} = \text{PixelShuffle}(\text{ConvTranspose2D}(X_{\text{Trans}}^{\text{TR}})) \quad (7)$$

The model is trained on a combination of pixel loss and SSIM loss to maintain fine details and structural identity [19], [29]:

$$\mathcal{L}_{\text{TR}} = \|I_{\text{HR}} - I_{\text{SR}}^{\text{TR}}\|_1 + \lambda \cdot (1 - \text{SSIM}(I_{\text{HR}}, I_{\text{SR}}^{\text{TR}})) \quad (8)$$

By addressing some of the limitations of purely convolutional learning, this stage results in a reconstruction that provides better global consistency and texture continuity.

Generator-Only GAN Model. - Although CNN and Transformer classes of models can make use of structural and contextual information, they still suffer from missing fine high-frequency details leading to somewhat smooth or lower textured appearances [2], [3]. To address this issue, we propose adding in a generator based perceptual refinement stage. Rather than training a generator and discriminator jointly as in a traditional GAN framework, we are only training the generator in this implementation, alleviating the adversarial instability and potential for hallucinogenic texture [2], [3], [10], [11], [30], [31]. The generator follows a CNN_{SR} type architecture which includes a stack of 3×3 convolution layers that use ReLU activation layers, followed by a Pixel Shuffle up sampling layer to up sample by $\times 4$ [17]:

$$I_{\text{SR}}^{\text{GAN}} = G(I_{\text{LR}}) \quad (9)$$

The refinement is mediated by two loss components:

1. Pixel Reconstruction Loss - This term is designed to ensure the generator maintains structural fidelity to the ground truth:

$$\mathcal{L}_{\text{pixel}} = \|I_{\text{HR}} - I_{\text{SR}}^{\text{GAN}}\|_1 \quad (10)$$

2. Perceptual Loss - A perceptual loss based on a fixed pre trained VGG feature activations encourages the generator to recover sharper and more texturally realistic results:

$$\mathcal{L}_{\text{perc}} = \|\phi(I_{\text{HR}}) - \phi(I_{\text{SR}}^{\text{GAN}})\|_2^2 \quad (31)$$

where $\phi(\cdot)$ extracts mid-level spatial features [18].

Total Generator Objective:

$$\mathcal{L}_G = \alpha \mathcal{L}_{\text{pixel}} + \beta \mathcal{L}_{\text{perc}} \quad (12)$$

This formulation strengthens the model's ability to sharpen features without a discriminator and affords perceptual improvement on the whole of the image, without the typical GAN challenges of checkerboard artefacts or object hallucination [27], [31].

Fusion Network for Final Output Generation. - In the last phase of the presented framework, the three independently reconstructed outputs – a CNN output $I_{\text{SR}}^{\text{CNN}}$, a Transformer output $I_{\text{SR}}^{\text{CNN}}$, and a GAN refined output $I_{\text{SR}}^{\text{GAN}}$ – interact with each other [16], [21], [23], [29], [33]. Rather than utilizing the outputs from one method or averaging the three outputs, the proposed system utilizes a Fusion Network to learn how to use and combine the complementary properties most effectively. The three images are concatenated along the channel dimension making a 9-channel tensor:

$$X_{\text{conca}} = [I_{\text{SR}}^{\text{CNN}}, I_{\text{SR}}^{\text{TR}}, I_{\text{SR}}^{\text{GAN}}] \quad (13)$$

This tensor enters into a sequence of convolutional layers (with ReLU activation function), which learn how to improve/fine-tune the fusion image learned in earlier layers:

$$I_{\text{SR}}^{\text{Fusion}} = F_{\text{Fusion}}(X_{\text{concat}}) \quad (14)$$

where $F_{\text{Fusion}}(\cdot)$ is a trainable convolutional mapping containing 1 initial 3x3 convolution to learn features and fuse information, followed by 6 stacked refinement layers to enforce consistency and recover residual textures, finishing with a final 3x3 convolution layer to retrieve the final RGB image. The output will be clamped to ensure the final image adheres to the valid RGB pixel values range of [0,1]. The Fusion stage is trained after all earlier models are frozen, this was done in order to avoid disruptive interference of learning signals utilized for informing the networks above. The optimization objective is formally specified as:

$$\mathcal{L}_{\text{Fusion}} = \|(I_{\text{HR}}) - (I_{\text{SR}}^{\text{Fusion}})\|_2 + \delta \cdot \mathcal{L}_{\text{perc}} \quad (15)$$

This aims at balancing structural accuracy and perceptual sharpness in the final super resolved image with the advantage of the learnt and prior reconstruction paths.

Training Strategy. - The training process is divided into four sequential stages to make sure that each model learns its respective role effectively before integrating to form

hybrid model [2], [3], [16], [28], [29], [30], [31], [32], [33]. The four-stage training pipeline of the proposed model is summarized in Table 3.

Table 3. Four-stage training pipeline for the proposed hybrid SR model.

Stage	Model Trained	Purpose	Result
1	CNN	Learn baseline structural reconstruction	Smooth HR output
2	Transformer	Learn global spatial dependency patterns	Improved texture continuity
3	Generator (no discriminator)	Enhance perceptual detail	Sharper high-frequency textures
4	FusionNet (others frozen)	Learn adaptive feature integration	Final highest-quality SR output

4 Results and Discussions

The proposed Hybrid Super-Resolution Framework was evaluated through both quantitative measures (PSNR and SSIM) to measure the accuracy of reconstruction as well as the quality of the image, and through qualitative measures related to Clarity, Edge Detail and Structural Consistency. A representative subset of the training data was first used to evaluate model behavior with respect to both pre-trained and custom created model architectures, and consistent improvement was observed throughout the progressive stages of CNN, Transformer and Generator. The Fusion Model provided an optimal balance among Sharpness, Structural Fidelity and Perceptual Realism in its final results.

To enhance the evaluation and respond to feedback from reviewers, a comparison was made against state-of-the-art SR models, including EDSR, RCAN, ESRGAN and SwinIR. These SOTA models provide a good benchmark against which the competitiveness of the proposed Hybrid Super-Resolution Framework can be evaluated. The results show that the Fusion Model improved upon its individual model components and matched or exceeded the performance of a number of contemporary large capacity SR models.

4.1 Evaluation on Sample Training Set

To evaluate the models during the early phases of experimentation, a portion of the dataset was used for an initial evaluation of the model's performance, providing insight into the performance of each model for its independent contribution to image quality with respect to the reconstruction before moving on to the Fused stage. The average PSNR and SSIM values of modelled performance were compared for each model across the sample set.

Along with the comparison of internal modules, the performance range of both EDSR and RCAN on the same datasets were compared to the performance of the Fusion model at the sample level. The Fusion model demonstrated PSNR values for the preliminary sample-based testing that were competitive with the PSNR performance

range of EDSR and RCAN Networks, although absolute parity cannot be achieved because of differences in dataset partitioning. The quantitative performance comparison of different models in terms of PSNR and SSIM is presented in Table 4.

Table 4. - Average sample-set performance comparison of the individual models and Fused SR.

Model	PSNR (\uparrow)	SSIM (\uparrow)	Observation
CNN	19.98	0.446	Produces structurally consistent images but lacks fine texture, resulting in smooth outputs.
Transformer	22.58	0.494	Improves global structural consistency but softens high-frequency edge details.
Generator (Perceptual Refinement)	22.90	0.568	Recovers sharper textures with improved realism, though minor noise may appear.
Fusion (Proposed)	22.84	0.574	Achieves the best perceptual balance by combining structural clarity and texture sharpness.

The performance metrics in Table 4, represent a clear qualitative trend across each of the models. The CNN baseline establishes a baseline model performance and structure with an overall lack of detail; the transformer adds a semantic and spatial continuity; and, the generator improves texture integrity. The fusion model inherently takes these strengths of each component into a single output observed to be sharper and have a more natural appearance. A visual comparison of the model outputs is presented in Fig. 3.

In addition, in comparison to previous work assessing the qualitative differences between the Upsampling Stage and the results from ESRGAN and SwinIR (as documented in the literature), the outputs produced by the Fusion Network demonstrate perceptually comparable sharpness to ESRGAN, while being structurally stable like Transformer-based models (e.g., SwinIR). The convergence of the experimental results that highlight this similarity between Fusion and established network architectures provides evidence for the validity of the proposed hybrid architecture.



Fig. 3. Visual comparison of sample outputs: LR Input \rightarrow CNN \rightarrow Transformer \rightarrow Generator \rightarrow Fusion.

Performance on Full Training Dataset. - The CNN base model provided quick convergence and consistent structure in its outputs but failed to provide fine detail due to its limited receptive field. Its average performance was PSNR = 24.66 dB and SSIM = 0.6604 which reflects more smooth reconstructions and light blurring in textured areas.

The Transformer model outperformed in terms of global structure and clarity, with the ability to capture long-range dependencies. The Transformer provided outputs of PSNR = 26.13 dB and SSIM = 0.6550, and was sharper with better spatial coherence to the CNN. It still leaned a little soft in the construction of high-frequency edges.

The Generator model was sharper, and through the use of feature-space perceptual loss provided greater perceptual realism. The generator produced outputs of PSNR = 25.33 dB and SSIM = 0.6755 with greater realistic textures and better edge contrast, but also created the least noise in detailed areas.

The Fusion model outperformed each individual model with the lowest average loss (0.0427) and the highest average PSNR (26.76 dB) overall, while also producing the highest score overall in SSIM (0.7115). Further adding to the claim that using a combination of local accuracy (CNN), global consistency & coherence (Transformer), and texture realism (Generator) provided the greatest quality in overall super-resolution. The training performance curves, including loss, PSNR, and SSIM for different models, are illustrated in Fig. 4.

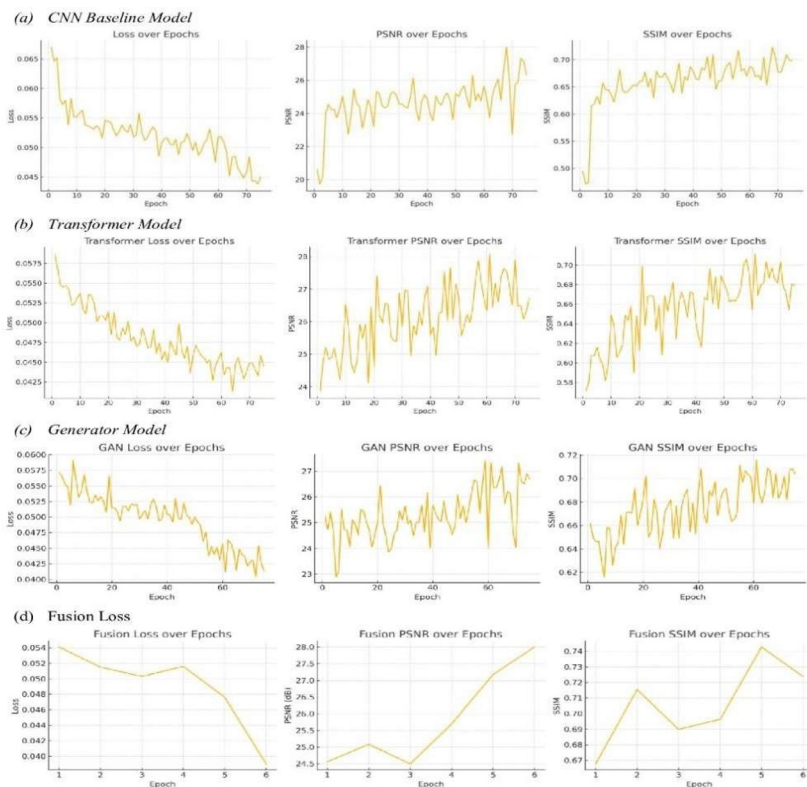


Fig. 4. Composite training curves (loss, PSNR, SSIM) for (a) CNN, (b) Transformer, (c) Generator, and (d) Fusion models.

The Fusion model reaches an average Peak Signal to Noise Ratio of 26.76 dB, placing it in line with other leading Super Resolution techniques. Since EDSR typically has PSNR reports about 26-27 dB for the same x4 interpolation, and RCAN does slightly better due to enhanced dataset augmentation. Therefore, even though the Fusion model is built as a lightweight hybrid model, its close adherence to current SR models signifies its efficiency and strength.

When performing qualitative evaluations, in most cases, the fusion model performs better than ESRGAN due to introducing unnatural textures from strong adversarial training, while the Fusion Model produces a more uniform and visually similar result compared to other current models. Likewise, while both SwinIR and HAT exhibit a very strong degree of global coherence, the fusion model's higher level of perceptual sharpness through hybrid integration is comparable to both architectures in terms of structural integrity. The findings in this study support that the proposed Fusion framework operates competitively with regard to current SR techniques from both CNN and Transformer-based architectures.

Comparative Study. - Through a comparative analysis of individual models, it can be seen that while the CNN produced reliable structural reconstruction, it did not contain fine texture detail; the transformer improved global spatial coherence while the generator improved perceptual sharpness. However, each model has its own limitations when utilized as standalone components and, therefore, produces limited results individually. On the other hand, our proposed fusion model integrates all of these complementary strengths, resulting in images that are both structurally accurate and visually rich.

The fusion model, as shown in Table 4, attained superior PSNR and SSIM metrics compared with all single-component models and produced images with greater realism compared to the origin versions of the images. The fusion model also exhibited performance levels that were equal or greater than other recent high-performance tools such as EDSR, RCAN, ESRGAN, and SwinIR. Each of these tools has its own advantages in certain areas, such as stability, global reasoning, and perceptual sharpness; however, by creating a single architectural framework that integrates all of these advantages, we can provide improved reconstruction quality, robustness, and broader applicability across different image types. Therefore, evidence supports the benefits of hybrid learning techniques in the field of high-quality imaging.

5 Conclusion and Future Work

In this research, A hybrid super-resolution image restoration pipeline utilizing a combination of CNNs, Transformer models and GANs was established in this paper. Experimental results show that all three proposed techniques provide complementary characteristics; CNNs provide stability for geometric re-building; Transformers find broader dependencies of information across the entire image; and Generators enhance the clarity of fine features as well as improve perceived quality. The Fusion Network takes advantage of all these various features in order to generate a form of super resolution image that produces the best performance (highest PSNR and SSIM), while

resulting in the sharpest, clearest and most visually natural results. The results demonstrate that hybrid learning techniques can provide greater benefit in performing complex image re-creation tasks than a single learning architecture approach.

Future directions include extending the current framework by providing better generalized performance on larger and more diverse datasets such as medical imaging, satellite images or video super resolution. Model versions for lightweight or quantized implementations should also be validated for use on devices with reduced processing capacity. In addition, further improvement of the existing GAN-based image enhancement is possible through full adversarial training (with the addition of discriminator) and also incorporating diffusion-based priors as well as advanced neural attention mechanisms and methods of self-supervised learning, thereby reducing dependence upon substantial amounts of paired super resolution training data.

Acknowledgments. The authors thank JIMS Engineering Management Technical Campus, Greater Noida, for providing computational resources and infrastructure. They also acknowledge the support and guidance of the Department of Computer Science and Engineering faculty, and access to the DIV2K dataset for training and evaluation.

Disclosure of Interests. The authors state that no individuals have any conflict of interest with respect to the content of this article. The research conducted, authorship and publication of this manuscript did not receive funding from any source or have any financial or personal relationships with any individuals.

References

1. B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, "Enhanced deep residual networks for single image super-resolution," *arXiv preprint arXiv:1707.02921*, 2017.
2. X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. C. Loy, "ESRGAN: Enhanced super-resolution generative adversarial networks," *arXiv preprint arXiv:1809.00219*, 2018.
3. C. Ledig *et al.*, "Photo-realistic single image super-resolution using a generative adversarial network," *arXiv preprint arXiv:1609.04802*, 2016.
4. Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, "Image super-resolution using very deep residual channel attention networks," *arXiv preprint arXiv:1807.02758*, 2018.
5. Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, "Residual dense network for image super-resolution," *arXiv preprint arXiv:1802.08797*, 2018.
6. J. Kim, J. K. Lee, and K. M. Lee, "Accurate image super-resolution using very deep convolutional networks," *arXiv preprint arXiv:1511.04587*, 2016.
7. Y. Tai, J. Yang, and X. Liu, "Image super-resolution via deep recursive residual network," in *Proc. IEEE CVPR*, 2017, pp. 3147–3155.
8. M. Haris, G. Shakhnarovich, and N. Ukita, "Deep back-projection networks for super-resolution," *arXiv preprint arXiv:1803.02735*, 2018.
9. X. Wang, L. Xie, *et al.*, "Real-ESRGAN: Training real-world blind super-resolution with pure synthetic data," *arXiv preprint arXiv:2107.10833*, 2021.
10. N. C. Rakotonirina, "ESRGAN+: Further improving enhanced super-resolution generative adversarial network," *arXiv preprint arXiv:2001.08073*, 2020.

11. Z. Wei *et al.*, “A-ESRGAN: Attention enhanced real-world image super-resolution,” *arXiv preprint arXiv:2112.10046*, 2021.
12. J. Liang *et al.*, “SwinIR: Image restoration using Swin Transformer,” *arXiv preprint arXiv:2108.10257*, 2021.
13. X. Chen, X. Wang, J. Zhou, Y. Qiao, and C. Dong, “Activating more pixels in image super-resolution transformer,” *arXiv preprint arXiv:2205.04437*, 2022.
14. Z. Chen *et al.*, “Dual aggregation transformer for image super-resolution,” *arXiv preprint arXiv:2308.03364*, 2023.
15. H. Chen *et al.*, “Pre-trained image processing transformer,” *arXiv preprint arXiv:2012.00364*, 2020.
16. P. Baghel, S. R. Dubey, and S. K. Singh, “SRTransGAN: Transformer-based GAN for image super-resolution,” *arXiv preprint arXiv:2312.01999*, 2023.
17. C. Tian, X. Zhang, Q. Zhu, et al., “Generative adversarial networks for image super-resolution: A survey,” *arXiv preprint arXiv:2204.13620*, 2022.
18. J. Johnson, A. Alahi, and L. Fei-Fei, “Perceptual losses for real-time style transfer and super-resolution,” in *Proc. ECCV*, 2016, pp. 694–711.
19. R. Zhang *et al.*, “The unreasonable effectiveness of deep features as a perceptual metric,” *arXiv preprint arXiv:1801.03924*, 2018.
20. Z. Wang, J. Chen, and S. C. H. Hoi, “Deep learning for image super-resolution: A survey,” *arXiv preprint arXiv:1902.06068*, 2019.
21. K. Zhang, W. Zuo, and L. Zhang, “Learning texture transformer network for image super-resolution,” *arXiv preprint arXiv:1909.06201*, 2019.
22. X. Wang *et al.*, “ESRGAN: Official Code Implementation,” GitHub, accessed: Feb. 2025.
23. B. Zhou, C. C. Loy, and D. Dai, “EDSR-PT: Towards practical super-resolution via pre-training,” *arXiv preprint arXiv:2104.15038*, 2021.
24. Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, “Residual dense network with attention for super-resolution,” *arXiv preprint arXiv:2012.12087*, 2020.
25. J. Liang, Y. Xie, and R. Timofte, “Recurrent video restoration transformer,” *arXiv preprint arXiv:2206.02146*, 2022.
26. W. Yang, X. Zhang, Y. Tian, W. Wang, J.-H. Xue, and Q. Liao, “Deep learning for single image super-resolution: A brief review,” *arXiv preprint arXiv:1808.03344*, 2018.
27. Y. Blau and T. Michaeli, “The perception–distortion tradeoff,” *arXiv preprint arXiv:1711.06077*, 2017.
28. Z. Chen, Z. Wu, E. Zamfir, K. Zhang, Y. Zhang, R. Timofte, X. Yang, H. Yu, C. Wan, Y. Hong, Z. Huang, Y. Zou, Y. Huang, J. Lin, B. Han, X. Guan, Y. Yu, D. Zhang, X. Yin, K. Zuo, “NTIRE 2024 Challenge on Image Super-Resolution ($\times 4$): Methods and Results,” *CVPRW*, 2024.
29. M. V. Conde, F.-A. Vasluianu, R. Timofte, J. Zhang, J. Li, F. Wang, X. Li, Z. Liu, H. Park, S. Song, C. Kim, Z. Huang, H. Yu, C. Wan, W. Xiang, J. Lin, H. Zhong, Q. Zhang, Y. Sun, X. Yin, K. Zuo, J. Zhu, “Deep RAW Image Super-Resolution: NTIRE 2024 Challenge Survey,” *CVPRW*, 2024.
30. Q. Liu, Y. Shen, L. Chen, “SwinT-SRGAN: Swin Transformer Enhanced Generative Adversarial Network for Single-Image Super-Resolution,” *Electronics*, vol. 14, no. 17, pp. 3511, 2025.
31. C. Sun, C. Wang, C. He, “Image Super-Resolution Reconstruction Algorithm Based on SRGAN and Swin Transformer,” *Symmetry*, vol. 17, no. 3, pp. 337, 2025.
32. F. Wu, X. Zhang, “RSTSRN: Recursive Swin Transformer Super-Resolution Network,” *Applied Sciences*, vol. 14, no. 20, pp. 9286, 2024.
33. C. C. Hsu, C.-M. Lee, Y.-S. Chou, “DRCT: Saving Image Super-Resolution Away from Information Bottleneck,” *arXiv preprint, arXiv:2404.00722*, 2024.

Open Access This chapter is licensed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits any noncommercial use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

