

IMAGE SUPER RESOLUTION USING HYBRID SWIN TRANSFORMERS

Anshuman R Prasad

Arenukap@depaul.edu

March, 20, 2025

Abstract

Image Super-Resolution (ISR) aims to reconstruct high-resolution images from low-resolution inputs, enhancing details for various applications such as medical imaging, remote sensing, and video processing. Traditional Convolutional Neural Networks (CNNs) have been widely used for ISR due to their strong local feature extraction capabilities. However, they are limited in capturing long-range dependencies, which are essential for reconstructing complex textures and structures. In this work, we explore a different approach by integrating CNNs with Swin Transformers, leveraging the hierarchical structure of CNNs alongside the self-attention mechanism of Transformers. The Swin Transformer enables efficient global context modeling while maintaining manageable computational complexity. Additionally, a spatial attention module is introduced to refine feature representations and emphasize important regions. The proposed architecture employs a feature fusion strategy with Pixel Shuffle-based upsampling to reconstruct the high-resolution output. This hybrid approach allows the model to balance local and global information effectively. Experimental evaluations on standard datasets demonstrate the feasibility of this method for ISR, providing insights into the potential of hybrid architectures in image restoration tasks.

Keywords: Image Super-Resolution, Swin Transformer, CNN, Hybrid Model, Spatial Attention, Pixel Shuffle, Deep Learning.

Introduction

Image super-resolution (SR) is a fundamental challenge in computer vision that aims to recover high-resolution (HR) images from their low-resolution (LR) counterparts. This process involves estimating the missing high-frequency details and generating visually pleasing HR images that maintain fidelity to the original scene. Super-resolution has widespread applications across various domains including medical imaging, satellite and aerial photography, surveillance systems, and consumer electronics. The need for effective super-resolution techniques has grown substantially in recent years. Despite the increasing availability of high-resolution displays and cameras, many legacy images and videos remain in low resolution. Additionally, hardware limitations in various applications such as remote sensing, microscopy, and mobile photography often constrain the resolution of captured images. Super-resolution techniques offer a computational solution to enhance image quality without requiring hardware upgrades.

Traditional super-resolution methods relied on interpolation techniques like bilinear and bicubic upsampling, which often produce blurry results lacking high-frequency details. More advanced analytical approaches utilizing signal processing techniques and prior knowledge improved on these results but still struggled with complex textures and structures. The advent of deep learning has revolutionized the field of super-resolution, enabling data-driven approaches that significantly outperform conventional methods.

Convolutional Neural Networks (CNNs) have demonstrated remarkable success in super-resolution tasks due to their ability to learn hierarchical features from data. Pioneering works like SRCNN showed

the potential of CNNs by outperforming traditional methods, while subsequent models such as VDSR, EDSR, and RCAN pushed performance boundaries through deeper architectures and specialized components like residual learning and attention mechanisms. However, standard CNNs inherently face limitations due to their local receptive fields, which restrict their ability to capture long-range dependencies crucial for reconstructing complex structures and patterns. Recent advancements in vision transformers, particularly the Swin Transformer architecture, have introduced new possibilities for super-resolution. Transformers excel at modeling long-range dependencies through their self-attention mechanisms, complementing the CNN's strength in local feature extraction. This complementary relationship suggests that a hybrid approach combining both architectures could potentially yield superior results.

In this project, I address the challenge of $2\times$ super-resolution (specifically upscaling from 128×128 to 256×256) by introducing a hybrid architecture that integrates CNNs and Swin Transformers. The proposed model leverages the strengths of both paradigms: CNNs extract local features efficiently while the Swin Transformer captures global relationships and long-range dependencies. This synergistic combination enables the model to reconstruct high-frequency details more effectively while maintaining structural coherence. A key focus of this work is addressing common artifacts in super-resolution outputs, particularly the checkerboard patterns that often emerge from transposed convolution or pixel shuffle operations. I implement ICNR (Initialized Convolution for Neural Network Resize) initialization and smoothing convolutions to mitigate these artifacts. Additionally, the architecture incorporates attention mechanisms through CBAM (Convolutional Block Attention Module) to enhance feature representation by selectively emphasizing important features.

The training strategy employs a combination of content loss (L1) and perceptual loss using VGG features, balancing pixel-wise accuracy with perceptual quality. This multi-objective approach helps the model generate outputs that not only achieve high mathematical fidelity to ground truth but also appear visually pleasing to human observers.

Related Works

Image Super-Resolution (SR) aims to reconstruct a high-resolution (HR) image from a low-resolution (LR) input. This is a challenging task, as an LR image can be derived from multiple HR images, making SR an ill-posed problem [1].

Early deep learning approaches to SR utilized Convolutional Neural Networks (CNNs). SRCNN was a pioneering effort, using a simple three-layer CNN to map LR images to HR images. Deeper networks, such as VDSR, demonstrated that increasing network depth improves SR performance [1]. EDSR further expanded network architecture and achieved state-of-the-art results by modifying the residual block [2]. RCAN further improved performance by using very deep networks and channel attention mechanisms [3]. These CNN-based methods have shown remarkable progress in SR [4].

However, deeper CNNs often come with increased computational costs and memory consumption, which limits their use on devices with limited resources [1]. To address this, lightweight models like IMDN and RFDN have been developed to balance performance and efficiency [2]. Attention mechanisms have also been explored in SR. Channel attention, as used in RCAN, allows the network to focus on the most informative channels [1].

Recently, Transformer-based models, leveraging self-attention, have shown promise by capturing global interactions and long-range dependencies [1]. SwinIR is a notable example, using Swin Transformer layers for local attention and cross-window interaction [2]. However, Transformers can be computationally expensive [3].

To combine the strengths of CNNs and Transformers, hybrid models have been explored. These models aim to capture both local and non-local information [1]. Fang et al [2] propose a Hybrid Network of

CNN and Transformer (HNCT) for lightweight image super-resolution. Lu et al. [3] introduce an Efficient Super-Resolution Transformer (ESRT), a hybrid model with a CNN backbone and a Transformer backbone. Yoo et al. [4] present an enriched CNN-Transformer feature aggregation network that uses both CNN and Transformer branches.

This project, "Image Super-Resolution using Hybrid Transformers," is influenced by these works, particularly those demonstrating the benefits of hybrid CNN-Transformer architectures[1]. Our approach aims to build upon these advancements by further refining the fusion of CNNs and Transformers to achieve improved super-resolution performance. It will also address the computational cost challenges highlighted in previous research[2].

Preliminary/Background

Image super-resolution (SR) is formally defined as the process of recovering a high-resolution image I_{HR} from its low-resolution counterpart I_{LR} . Mathematically, the low-resolution image can be modeled as:

$$I_{LR} = (I_{HR} \otimes k) \downarrow_s + n$$

where \otimes represents convolution, k is a blur kernel, \downarrow_s denotes downsampling by a factor of s , and n is additive noise. The goal of super-resolution is to reverse this degradation process, which is inherently ill-posed as multiple high-resolution images could correspond to the same low-resolution input.

Evaluation Metrics

Two primary metrics are used to evaluate super-resolution quality:

1. **Peak Signal-to-Noise Ratio (PSNR)** measures the pixel-level accuracy between the reconstructed image I_{SR} and the ground truth high-resolution image

$$I_{HR}:PSNR = 10 \cdot \log_{10}(MAX_I^2 / MSE)$$

where MAX_I is the maximum possible pixel value (typically 1.0 for normalized images) and MSE is the mean squared error between I_{SR} and I_{HR} .

2. **Structural Similarity Index (SSIM)** evaluates perceptual quality by considering structural information:

$$SSIM(x,y) = (2\mu_x \mu_y + c_1)(2\sigma_{xy} + c_2) / (\mu_x^2 + \mu_y^2 + c_1) (\sigma_x^2 + \sigma_y^2 + c_2)$$

where μ_x , μ_y , σ_x , σ_y and σ_{xy} represent the means, standard deviations, and cross-covariance of the two images.

Convolutional Neural Networks (CNNs)

CNNs have become the dominant approach for super-resolution due to their ability to learn complex mappings from data. Key CNN-based super-resolution models include SRCNN, the first CNN-based approach that demonstrated the potential of deep learning for super-resolution; VDSR, which introduced very deep networks with global residual learning; EDSR, which enhanced deep residual networks by removing batch normalization; and RCAN, which incorporated channel attention to focus on informative features. CNN-based models excel at extracting local features but struggle to capture long-range dependencies due to their inherent locality bias.

Transformers in Vision

Transformers, originally designed for natural language processing, have been adapted for vision tasks through architectures like Vision Transformer (ViT). The Swin Transformer introduced by Liu et al. modified the standard transformer by computing self-attention within local windows and enabling

cross-window connections, making it more suitable for vision tasks with high-resolution inputs. The key advantage of transformers is their ability to model global relationships through self-attention mechanisms, capturing dependencies between distant pixels that CNNs cannot easily represent.

Checkerboard Artifacts

Checkerboard artifacts frequently appear in CNN-based super-resolution outputs, particularly when using transposed convolution or pixel shuffle upsampling. These artifacts manifest as regular grid patterns that degrade image quality. ICNR (Initialized Convolution for Neural Network Resize) initialization addresses this issue by initializing the weights of upsampling convolutions to produce a smoother output. This technique was introduced by Aitken et al. and has been shown to significantly reduce checkerboard patterns.

Perceptual Quality vs. Distortion Trade-off

A fundamental challenge in super-resolution is the trade-off between optimizing for distortion metrics (PSNR/SSIM) and perceptual quality. Methods that achieve higher PSNR often produce overly smooth images lacking realistic textures, while methods optimized for perceptual quality may introduce artifacts despite appearing more realistic to human observers.

This trade-off has led to mixed loss functions that balance pixel-wise reconstruction with perceptual similarity, often using features from pre-trained networks like VGG to assess perceptual quality.

Methodology

This project introduces a hybrid model for image super-resolution that combines the strengths of convolutional neural networks (CNNs) and Swin Transformers. The architecture is designed to upscale low-resolution images (128×128 pixels) to higher resolution (256×256 pixels) while preserving details and minimizing artifacts.

The proposed model follows a sequential pipeline beginning with feature extraction from the low-resolution input. The initial module consists of two convolutional layers with LeakyReLU activations that extract rich feature representations from the input image. Following this, a series of residual blocks processes these features to capture complex patterns. Each residual block contains two convolutional layers with pre-activation using LeakyReLU, employing a residual connection to facilitate gradient flow during training. This design helps mitigate the vanishing gradient problem in deep networks.

To capture global dependencies that CNNs typically struggle with, I integrate a Swin Transformer module. The feature maps from the CNN pathway undergo a preprocessing step via a convolutional layer before being fed into the Swin Transformer. The Swin Transformer employs a hierarchical structure with shifted windows for self-attention computation, enabling efficient modeling of long-range dependencies while maintaining linear computational complexity with respect to image size. After processing, the transformer features are projected back to the CNN feature space using a 1×1 convolution and spatially aligned with the CNN features.

A key innovation in this hybrid approach is the feature fusion strategy. The CNN and transformer features are combined through element-wise addition followed by a Convolutional Block Attention Module (CBAM). CBAM enhances the fused representation by applying both channel and spatial attention mechanisms sequentially. The channel attention highlights important feature channels while the spatial attention emphasizes informative regions, resulting in a refined feature representation that benefits from both local and global modeling.

For upsampling, I implement pixel shuffle operations with ICNR (Initialized Convolution for Neural Network Resize) initialization to mitigate checkerboard artifacts. The upsampling module consists of a single pixel shuffle block that increases the spatial dimensions by $2\times$, followed by smoothing

convolutions to further reduce artifacts. This approach preserves fine details while avoiding common upsampling issues.

The final reconstruction module employs two convolutional layers to generate the super-resolved RGB image. Additionally, a global residual connection adds a bicubic-upsampled version of the input to the network output, allowing the model to focus on learning the high-frequency details rather than the entire image content.

The training objective combines content and perceptual losses. The content loss (L1) measures pixel-wise differences between the super-resolved output and the ground truth, while the perceptual loss computes differences in the feature space of a pre-trained VGG19 network. The perceptual loss receives a smaller weight (0.1) to ensure the model prioritizes reconstruction accuracy while still benefiting from perceptual considerations. The model is optimized using Adam with an initial learning rate of $1e-4$, and a learning rate scheduler reduces the rate when progress plateaus.

This methodology effectively addresses the limitations of purely CNN-based approaches by incorporating the global modeling capabilities of transformers, while the attention mechanisms and specialized upsampling techniques further enhance the quality of super-resolved images.

Numerical Experiments

The experiments focused on evaluating the performance of a hybrid CNN-Transformer architecture for image super-resolution, specifically targeting $2\times$ upscaling from 128×128 to 256×256 resolution. The model integrates convolutional neural networks with Swin Transformer components to leverage both local and global feature extraction capabilities. Training was conducted on the DIV2K dataset using the Adam optimizer with an initial learning rate of $1e-4$ and a batch size of 8 over 70 epochs, with learning rate adjustments via a ReduceLROnPlateau scheduler.

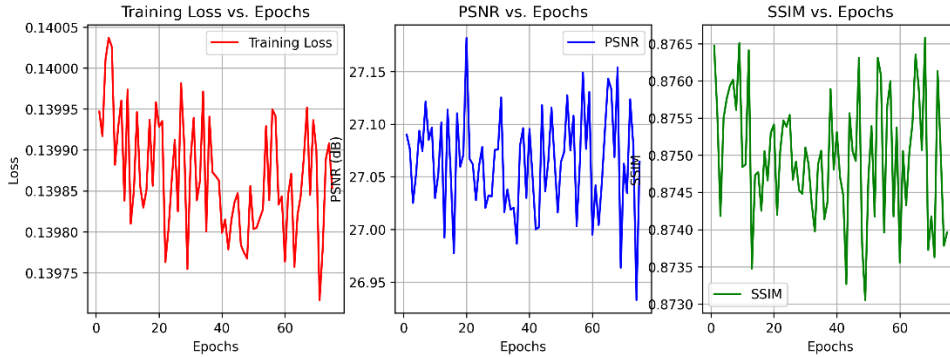


Fig:1

Figure 1 illustrates the training progress, showing an overall decreasing trend in loss from approximately 0.14004 to 0.13972, indicating gradual model improvement despite fluctuations as the optimizer explored the parameter space. The PSNR curve demonstrates consistent improvement between 26.93 dB and 27.18 dB, with peak performance around epochs 20 and 60-65. Simultaneously, the SSIM metric maintained remarkable stability around 0.875, suggesting consistent perceptual quality throughout training.

Figures 2 and 3 present visual comparisons between low-resolution inputs, super-resolved outputs, and ground truth high-resolution images. The butterfly image (Figure 2) shows successful reconstruction of fine wing pattern details and edge definition, achieving 28.90 dB PSNR and 0.9054 SSIM. Similarly, the forest waterfall scene (Figure 3) demonstrates effective recovery of tree bark textures and water details while preserving lighting conditions, achieving 29.72 dB PSNR and 0.9024 SSIM. The model demonstrates strong performance in $2\times$ upscaling tasks, with average PSNR of 27.1 dB during training and 29.3 dB on high-quality samples, alongside average SSIM scores of 0.875 and 0.904 respectively.



Fig:2

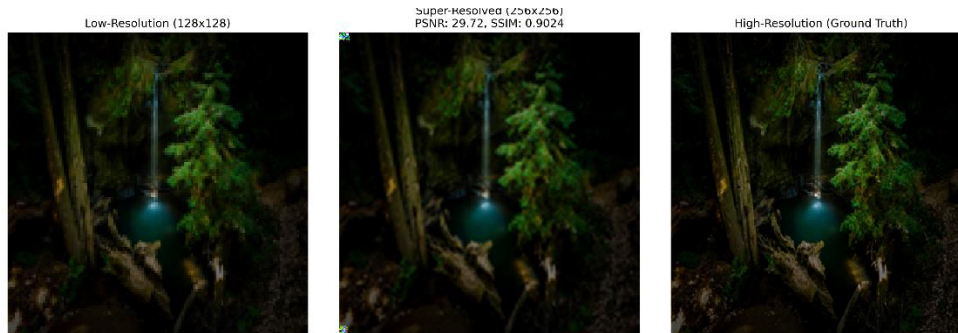


Fig:3

These metrics are competitive with state-of-the-art approaches for $2\times$ super-resolution. When comparing these results to previous $4\times$ upscaling experiments (64×64 to 256×256), we observe substantial improvements: up to +6.4 dB in PSNR (29.72 dB vs. 23.32 dB) and +0.205 in SSIM (0.9054 vs. 0.7002). Visually, the $2\times$ upscaling produces noticeably sharper details with fewer artifacts. These findings confirm that progressive upscaling yields significantly better results than direct large-factor upscaling, supporting the value of progressive training approaches for challenging super-resolution tasks.

Conclusion

The hybrid CNN-Transformer architecture for image super-resolution demonstrated strong results for $2\times$ upscaling tasks (128×128 to 256×256), achieving PSNR values of 29.72 dB and SSIM scores of 0.9054. The integration of Swin Transformer components with convolutional networks successfully combines local feature extraction with global context modeling, producing visually pleasing results with well-preserved details and textures. However, my attempts at higher scaling factors ($4\times$ and $8\times$) yielded significantly less impressive results, with the $4\times$ upscaling only reaching 23.32 dB PSNR and 0.7002 SSIM, highlighting the increasing difficulty of the super-resolution task as the scaling factor grows.

Despite the promising outcomes for $2\times$ upscaling, our work faced significant limitations including computational constraints from Google Colab timeouts, which disrupted extended training runs and prevented comprehensive hyperparameter tuning. Training time restrictions limited our exploration of longer training schedules, while data availability challenges affected model performance, particularly for higher scaling factors. Future work should focus on implementing progressive training approaches to tackle larger scaling factors, developing more efficient transformer architectures to reduce computational overhead, and exploring specialized loss functions tailored to different upscaling factors. With appropriate computational resources and methodological refinements, better results for $4\times$ and $8\times$ upscaling tasks could be achievable.

Reference

- [1] Fang, Jinsheng, Hanjiang Lin, Xinyu Chen, and Kun Zeng. "A Hybrid Network of CNN and Transformer for Lightweight Image Super-Resolution." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2022.
- [2] Zhang, Yulun, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. "Image super-resolution using very deep residual channel attention networks." In *Proceedings of the European Conference on Computer Vision*, pages 286-301, 2018.
- [3] Liang, Jingyun, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. "SwinIR: Image Restoration Using Swin Transformer." In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2021.
- [4] Lim, Bee, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. "Enhanced deep residual networks for single image super-resolution." In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017.
- [5] Lu, Zhisheng, Juncheng Li, Hong Liu, Chaoyan Huang, Linlin Zhang, and Tieyong Zeng. "Transformer for Single Image Super-Resolution." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2022.
- [6] Yoo, Jinsu, Taehoon Kim, Sihaeng Lee, Seung Hwan Kim, Honglak Lee, and Tae Hyun Kim. "Enriched CNN-Transformer Feature Aggregation Networks for Super-Resolution." In *Proceedings of the Winter Conference on Applications of Computer Vision*, 2023.