

Enhancing Drone Imagery with Super Resolution

A Final Project End Semester Report

submitted by

MARIYAM JOORY (CS20B1030)

in partial fulfilment of requirements

for the award of the degree of

BACHELOR OF TECHNOLOGY



**Department of Computer Science and Engineering
INDIAN INSTITUTE OF INFORMATION TECHNOLOGY,
DESIGN AND MANUFACTURING, KANCHEEPURAM**

May 2024

ACKNOWLEDGEMENTS

I would like to extend my sincerest gratitude to **Dr. Rahul Raman** for guiding me in a lot of decision making situations and giving valuable inputs to work on the project. His valuable insights and guidance were pivotal in shaping the project and enhancing my understanding in the field

During this time period I explored various models and metrics to develop and improve this project.

ABSTRACT

The field of single image super-resolution, focusing on reconstructing high-resolution images from low-resolution counterparts, is an actively advancing area within computer vision. This research delves into the latest advancements in super-resolution technologies, including models like Real-ESRGAN. A prevalent issue in these models is their limited ability to preserve the intricate textures and fine details critical for replicating real world scenes. This limitation often arises from their training on datasets primarily composed of less complex images, such as animations.

To overcome these challenges, this study conducts a comparative analysis of recent models to assess their capability in rendering true-to-life detail and texture. Furthermore, the research extends to practical tests using an enriched dataset of complex urban and natural scenes captured via drone photography. The goal is to refine the model's effectiveness for real-world applications, enhancing its utility in areas like environmental monitoring and urban development. The outcomes of this research are discussed, focusing on their implications for future studies and the enhancement of super-resolution techniques in realistic settings.

TABLE OF CONTENTS

ACKNOWLEDGEMENTS	i
ABSTRACT	ii
LIST OF TABLES	v
LIST OF FIGURES	vi
ABBREVIATIONS	vii
NOTATION	viii
1 Introduction	1
1.1 Background	1
1.2 Motivation	1
1.3 Objectives	2
2 Methodology	3
2.1 Models Employed for Enhancement	3
2.1.1 Bicubic	3
2.1.2 DeepCNN	3
2.1.3 Generative Adversarial Networks (GANs)	4
2.1.4 ESRGAN	4
2.1.5 Real-ESRGAN	5
2.1.6 BSRGAN	5
2.1.7 A-ESRGAN	6
2.1.8 SwinIR	6
2.2 Dataset Preparation and Selection	7
2.2.1 Quantitative Analysis	8

3 Work Done	11
3.1 Qualitative analysis	11
3.1.1 Detail and Texture Preservation	11
3.1.2 Artifact Reduction	11
3.1.3 Edge Sharpness	12
3.1.4 Colour Fidelity and Realism	12
3.2 Quantitative Analysis	15
3.2.1 Peak Signal To Noise Ratio	15
3.2.2 Structural Similarity Index	15
3.2.3 Learned Perceptual Image Patch Similarity (LPIPS)	15
3.3 Assessing New Models on Drone Images	16
3.4 Fine-Tuning Real-ESRGAN	19
3.4.1 Data Preparation	19
3.4.2 Hyperparameter Adjustments	19
3.4.3 Augmentation Techniques	19
3.5 Outcomes of Fine-Tuning Real-ESRGAN	22
3.5.1 Remaining Challenges	23
4 CONCLUSION AND FUTURE SCOPE	26
4.1 Conclusion	26
4.2 Future Scope	26
4.2.1 Exploring New Models and Techniques	26
4.2.2 Extended Applications	27
4.2.3 Enhancement of Model Training and Efficiency	27
4.2.4 Further Fine-Tuning and Optimization	27
REFERENCES	28

LIST OF TABLES

3.1	Comparison of Image Enhancement Models	16
3.2	Comparison of Super Resolution Models	18
3.3	Quantitative analysis of the two models	25

LIST OF FIGURES

2.1	Real-ESRGAN architecture	5
2.2	Swin-IR architecture	6
2.3	Campus photo example taken from a high vantage point.	7
2.4	Photos taken from around IIITDM	7
2.5	Photos taken from various datasets	8
3.1	Analysis of models on wall textures	12
3.2	Analysis of models on writing	13
3.3	Analysis of models on brick patterns	13
3.4	Analysis of models on natural foliage	14
3.5	Analysis of models on building structure	14
3.6	Analysis of models on window structures	17
3.7	Analysis of models on flooring	17
3.8	Analysis of models on letters	18
3.9	Analysis of models on fine details	18
3.10	Dataset for finetuning	20
3.11	Finetuning dataset with variable lighting conditions and elevations .	21
3.12	Analysis of models on building structure	23
3.13	Analysis of models on building structure	24
3.14	Analysis of models on distorted image	24
3.15	Analysis of models on distorted image	25
3.16	Analysis of models on building structure	25

ABBREVIATIONS

SSIM	Structural Similarity Index Measure
LPIPS	Learned Perceptual Image Patch Similarity
PSNR	Peak Signal to Noise Ratio
SRGAN	Super-Resolution Generative Adversarial Network
ESRGAN	Enhanced Super-Resolution Generative Adversarial Networks
CNN	Convolutional Neural Network
FSIM	Feature Similarity Index

NOTATION

$\phi_l(\cdot)$	Feature map at layer l of a pretrained neural network used in LPIPS.
w_l	Weight factor for layer l in the LPIPS metric.
μ_x, μ_y	Mean of pixels in windows x and y for SSIM.
σ_x^2, σ_y^2	Variance of pixels in windows x and y for SSIM.
σ_{xy}	Covariance of windows x and y for SSIM.
C_1, C_2	Stability constants in the SSIM calculation.

CHAPTER 1

Introduction

1.1 Background

The integration of unmanned aerial vehicles (UAVs), or drones, into diverse operational frameworks marks a significant shift from their initial recreational use to critical applications in several sectors. These applications have been enriched by advancements in artificial intelligence, autonomous navigation, and drone fleet management, which have significantly broadened the functional spectrum of UAV technology. Today, drones are pivotal in areas such as ecological monitoring, urban development, public safety, and defense due to their versatility and capability to operate under challenging conditions.

[1]

Drones, especially those enabled with AI, undertake complex tasks with greater autonomy, reducing the need for human oversight and risk exposure. Such systems are invaluable for monitoring natural habitats, inspecting infrastructure in remote areas, or surveillance operations sensitive zones. Additionally, the ongoing development of technologies for secure UAV communications, precise location tracking, and counter-drone strategies continually enhances the operational efficacy and safety of drone missions.

[2]

1.2 Motivation

One of the primary challenges in utilizing drone technology, especially in fields that require high precision, is often the sub optimal resolution of captured images due to various factors. This limitation can significantly diminish the utility and clarity of the data. There is a requirement for advanced super resolution technologies that can substantially upgrade image quality to meet the standards required by various professional

and scientific applications. Enhancing image resolution is important for tasks ranging from environmental conservation efforts to strategic urban planning and detailed security surveillance. [3]

1.3 Objectives

The objective of this study is to conduct a comprehensive assessment and refinement of different image enhancement models with special emphasis on super-resolution techniques designed to elevate the quality of drone-captured images. This research will evaluate various enhancement approaches to determine their efficacy in improving image detail, texture, and overall visual authenticity.

The objectives are laid out as follows:

1. Conduct an extensive comparative analysis of both established and emerging image enhancement models to determine their efficacy in enhancing image detail, texture, and realism.
2. Specifically fine-tune promising model to tailor their performance to the unique attributes of drone-captured images.
3. Develop a comprehensive framework for both qualitative and quantitative evaluation of these models to ensure the super-resolved images meet the necessary high-quality standards for their intended uses.

CHAPTER 2

Methodology

2.1 Models Employed for Enhancement

The chosen models encompass a range of techniques from traditional interpolation to cutting-edge deep learning approaches. This selection enables a comprehensive analysis of their effectiveness in enhancing image quality, with a particular focus on their application to real-world imagery commonly encountered in various terrains such as architectural structures and natural foliage. The comparative analysis of these models will contribute valuable insights into the field of image enhancement, informing best practices and guiding future research.

2.1.1 Bicubic

Bicubic interpolation is a resampling technique used for image scaling and enhancement. It works by utilizing the values of the nearest 16 pixels (4x4 environment) to estimate the new pixel value in a resized image. This method provides smoother transitions than nearest-neighbor or bilinear interpolation by considering the intensity values of surrounding pixels, resulting in higher quality and less jagged images. Bicubic interpolation is widely used due to its balance between computational efficiency and the quality of the output images, making it a standard choice in various image processing applications.

2.1.2 DeepCNN

Deep Convolutional Neural Networks (DeepCNNs) are a class of deep learning algorithms that excel in analyzing visual imagery. By leveraging multiple layers of processing, DeepCNNs can learn complex features at various levels of abstraction, from simple

edges to high-level objects within images. This ability makes DeepCNNs particularly effective for tasks like image recognition, classification, and enhancement, where understanding detailed visual information is crucial. Their deep architecture allows for the extraction of intricate patterns and details, enhancing image quality and resolution significantly. [4]

2.1.3 Generative Adversarial Networks (GANs)

Introduced in 2014 by Ian Goodfellow et al., Generative Adversarial Networks (GANs) are an innovative class of artificial intelligence algorithms. These networks comprise two components, the generator and the discriminator, which engage in a strategic interaction resembling a game theory scenario. The generator's objective is to fabricate data that is indistinguishable from genuine data, whereas the discriminator assesses whether the data comes from the generator or is authentic. The objective function for a GAN is defined as:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

where G denotes the generator, D the discriminator, x samples from the actual data distribution, z input noise samples for the generator, and V the value function that the discriminator strives to maximize and the generator seeks to minimize. This adversarial interaction propels both networks to continuously refine their methods until the generator can produce data that closely resembles real data. GANs have had a significant impact on the field of machine learning, providing essential tools for generating synthetic data, improving image quality, and more. [5]

2.1.4 ESRGAN

Enhanced Super-Resolution Generative Adversarial Networks (ESRGAN) marks a significant advancement in super-resolution technology. Utilizing a generative adversarial network (GAN) framework, ESRGAN comprises a generator that upscales images and a discriminator that assesses the quality of these enhanced images. The introduction

of residual-in-residual dense blocks (RRDB) within its architecture allows ESRGAN to capture and enhance fine textures with considerable detail, making it particularly suited for improving the resolution and quality of low-resolution images. [6]

2.1.5 Real-ESRGAN

Real-ESRGAN, an evolution of ESRGAN, focuses on achieving more realistic and natural-looking images by addressing common issues like artifacts and blurriness, further refining the quality of the upsampled images. Through architectural and training refinements, Real-ESRGAN enhances image quality by focusing on the authenticity and clarity of textures and edges, catering to applications that require high fidelity and realism in image enhancement. suited for improving the resolution and quality of low-resolution images. [7]

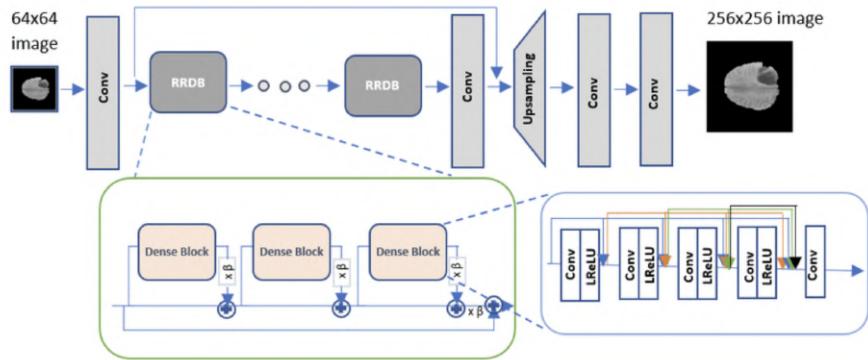


Figure 2.1: Real-ESRGAN architecture

2.1.6 BSRGAN

Blind Super Resolution Generative Adversarial Network (BSRGAN) is designed to tackle super resolution tasks under a variety of degradation models, unlike traditional methods that typically assume a single, known degradation. By training with a diverse set of degraded images, BSRGAN adapts to a wide range of quality issues in input images, making it robust in real-world scenarios where degradation is not uniform or known in advance. The network's ability to generalize from multiple types of degradations allows it to produce high-quality enhancements across a broader array of conditions compared to more traditional super-resolution methods.

2.1.7 A-ESRGAN

Attention Enhanced Super-Resolution Generative Adversarial Network (A-ESRGAN) incorporates attention mechanisms into the ESRGAN framework to focus more precisely on areas of an image that require intricate detailing during enhancement. This approach allows A-ESRGAN to manage resources more efficiently by prioritizing regions within an image that benefit the most from super-resolution techniques, such as areas with complex textures or significant degradation. The inclusion of attention mechanisms helps in maintaining texture fidelity and reducing artifacts, providing sharper and more detailed image outputs. [8]

2.1.8 SwinIR

SwinIR is a state-of-the-art image super resolution tool that leverages the Swin Transformer architecture, a novel approach that uses shifted windowing schemes to efficiently process image patches. This method brings the benefits of both global and local attention mechanisms, allowing SwinIR to capture contextual information over large areas while focusing on fine details within smaller regions. The flexibility and effectiveness of the Swin Transformer in handling various scales of details make SwinIR particularly effective for tasks requiring high precision in detail enhancement, such as medical imaging, satellite image interpretation, and advanced photographic restoration.[9]

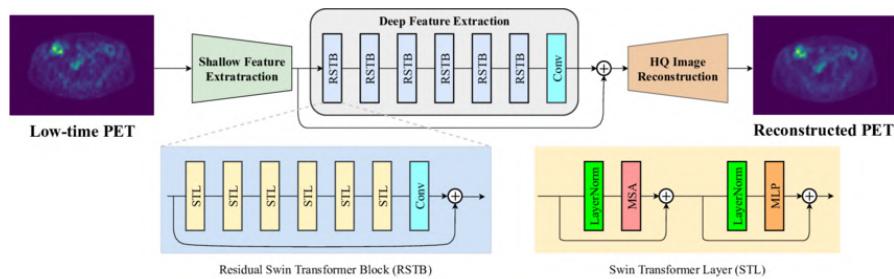


Figure 2.2: Swin-IR architecture

2.2 Dataset Preparation and Selection

For this study, we created a dataset 500 images using two methods for testing the models: photos taken around the campus and drone-shot images from the internet. The campus photos focus on buildings, trees, and various structures taken to capture the details and textures of these subjects under different lighting conditions. These images were taken from an high vantage points, to simulate perspective and range typically associated with drone photography.



Figure 2.3: Campus photo example taken from a high vantage point.

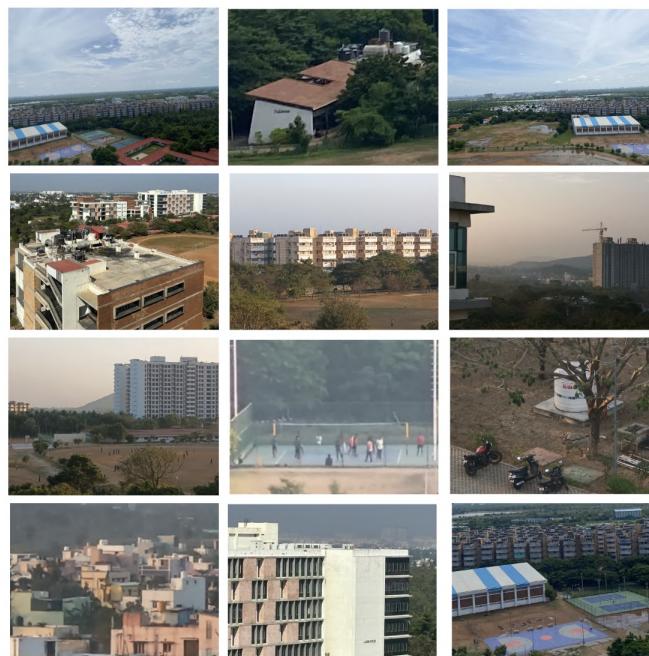


Figure 2.4: Photos taken from around IIITDM

The study utilized drone-shot image datasets sourced from several high-resolution aerial image collections designed for academic research, including the Stanford Drone Dataset [10]. These aerial datasets, complemented by ground-level photographs, encompass a diverse range of terrains including architectural settings and terrestrial. The combination of these datasets enables rigorous testing of enhancement models on images that present various challenges, such as expansive landscapes and differing resolutions typical of drone imagery.



Figure 2.5: Photos taken from various datasets

2.2.1 Quantitative Analysis

Peak Signal to Noise Ratio (PSNR)

PSNR is a critical metric in image processing that measures the relationship between the highest possible signal power and the power of corrupting noise impacting the image's fidelity. It evaluates the quality by comparing the original image, denoted as I , with its

enhanced counterpart, K . The PSNR is calculated using the formula:

$$\text{PSNR} = 10 \cdot \log_{10} \left(\frac{\text{MAX}_I^2}{\text{MSE}} \right)$$

where MAX_I represents the maximum potential pixel value within the image and MSE is the mean squared error between I and K . Generally, a higher PSNR value suggests a higher quality of the enhanced image.

Learned Perceptual Image Patch Similarity (LPIPS)

LPIPS assesses similarity by employing deep learning to approximate how the human eye perceives differences between two images. It calculates the perceptual distance between the original image I and its enhanced counterpart K :

$$LPIPS(I, K) = \sum_l w_l \cdot \frac{1}{H_l W_l} \sum_{h,w} |\phi_l(I)h, w - \phi_l(K)h, w|^2 \quad (2.1)$$

In this equation, $\phi_l(\cdot)$ denotes the feature maps from layer l of a trained neural network, and w_l represents the weighting factors for each layer.

Structural Similarity (SSIM) Index

The SSIM index is employed to assess the degree of similarity between two images by examining aspects such as structural integrity, brightness, and contrast. For two sections of the images, x from the original image I and y from the enhanced image K , each with dimensions $N \times N$, SSIM is computed as follows:

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)} \quad (2.2)$$

In this equation, μ_x and μ_y represent the mean values of pixels, σ_x^2 and σ_y^2 are the variances, and σ_{xy} is the covariance of x and y . Constants C_1 and C_2 are included to prevent division by a small denominator.

Feature Similarity Index (FSIM)

FSIM provides a metric that reflects critical visual quality aspects by incorporating factors like gradient magnitude and phase congruency, important components of human visual perception. The FSIM score between two images, I and K , is computed as follows:

$$FSIM(I, K) = \sum_{x,y} (S_L(x, y) \times PC_I(x, y)) \quad (2.3)$$

$S_L(x, y)$ evaluates the similarity based on local gradient and phase congruency between corresponding pixels, and $PC_I(x, y)$ signifies the phase congruency at pixel (x, y) of the original image. A higher FSIM score indicates a closer match to the human perception of image quality.

CHAPTER 3

Work Done

After setting up the models and preparing the dataset, we tested the four chosen image enhancement models—Bicubic, DeepCNN, ESRGAN, and Real-ESRGAN—on 500 images collected from different sources. These images covered a variety of subjects and conditions, making them suitable for evaluating how well each model improves image quality. The results of the model were run through quantitative and qualitative analysis.

3.1 Qualitative analysis

3.1.1 Detail and Texture Preservation

The Bicubic model's simplicity often leads to a loss of fine details, resulting in somewhat blurred textures. DeepCNN improves upon this by better preserving details, though it can render very fine textures too smoothly. ESRGAN enhances detail further but may introduce textures that diverge from the original image's authenticity. In contrast, Real-ESRGAN refines this approach, employing sophisticated techniques to enhance details more realistically and maintain the integrity of original textures.

3.1.2 Artifact Reduction

Bicubic interpolation is prone to pixelation and aliasing artifacts in detailed sections. DeepCNN, aiming to reduce artifacts, can inadvertently produce edge halos or ringing effects due to its focus on edges. ESRGAN generally diminishes artifact presence but isn't immune, particularly in textured areas. Real-ESRGAN excels by implementing advanced strategies that more effectively reduce artifacts especially JPEG compression artifacts, ensuring cleaner image enhancements.

3.1.3 Edge Sharpness

Bicubic often yields softer edges, lacking the precision for sharp delineation. Deep-CNN tends to overemphasize edges, sometimes resulting in an unnatural sharpness. ESRGAN offers an improvement with sharper edges, though it occasionally leads to edge artifacts. Real-ESRGAN provides a balanced solution, enhancing edges to be sharp and clear without exaggeration, ensuring a more natural appearance.

3.1.4 Colour Fidelity and Realism

Bicubic maintains color reasonably well but may introduce slight shifts. DeepCNN, despite enhancing overall image quality, can sometimes alter color balance or saturation. ESRGAN preserves colors more faithfully but may still cause minor shifts in hue or saturation. Real-ESRGAN stands out by focusing on color consistency, using advanced methods to ensure colors in enhanced images remain true to the original, contributing to overall realism.

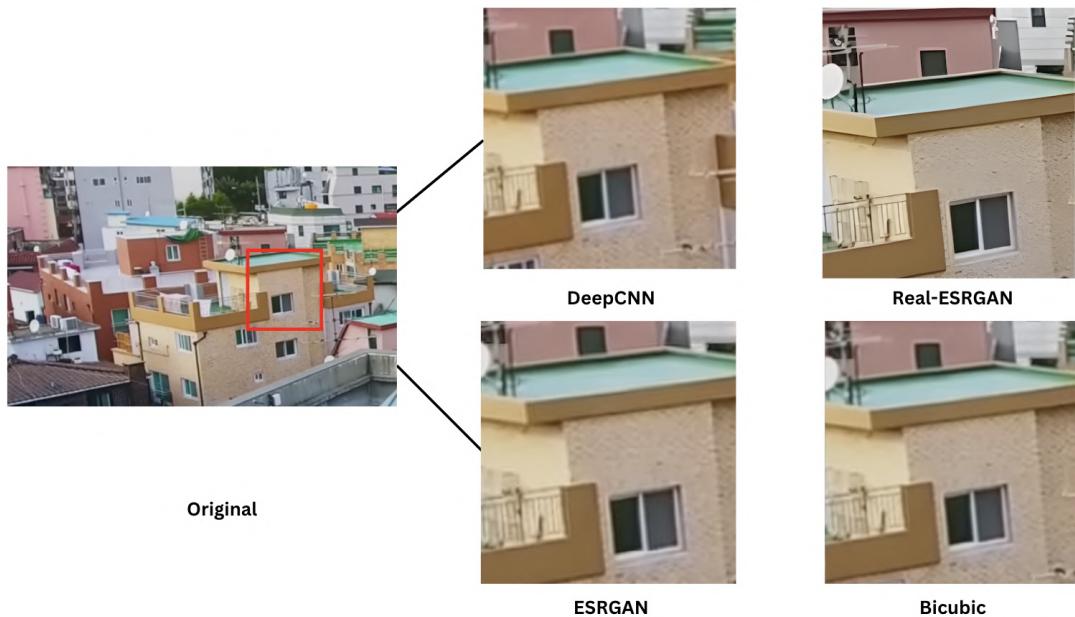


Figure 3.1: Analysis of models on wall textures

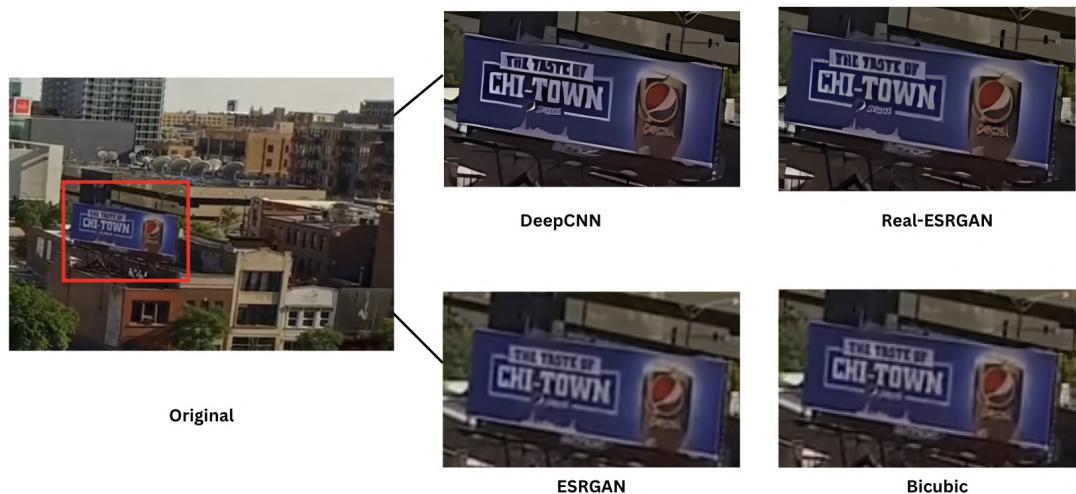


Figure 3.2: Analysis of models on writing

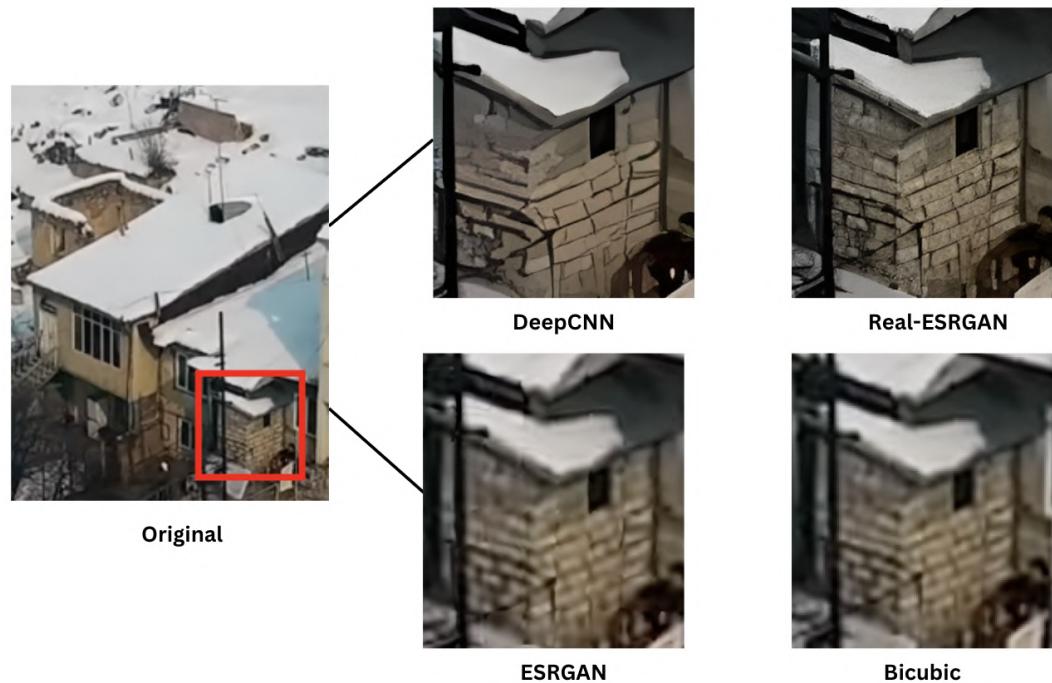


Figure 3.3: Analysis of models on brick patterns

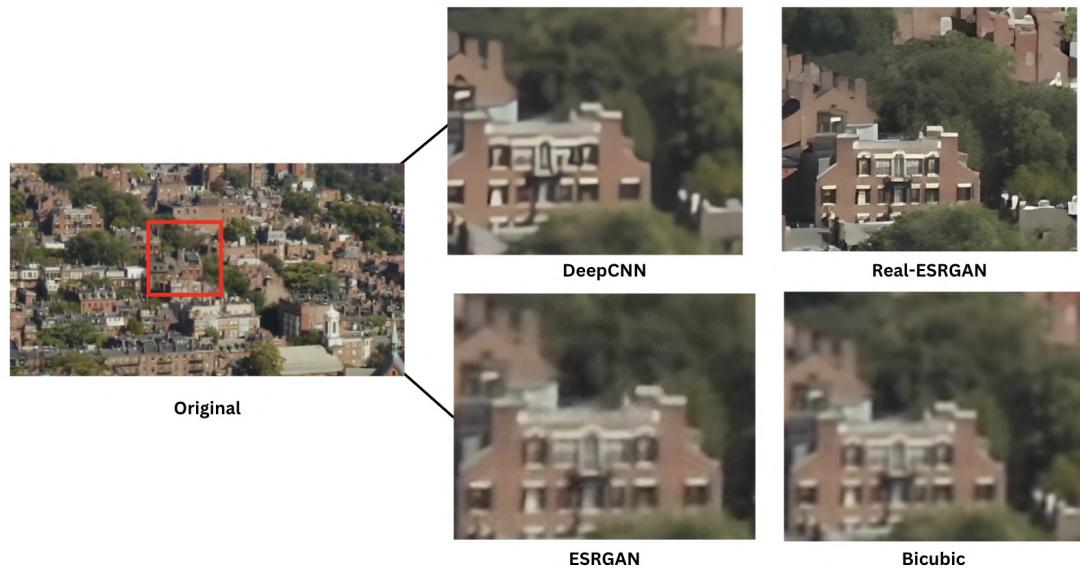


Figure 3.4: Analysis of models on natural foliage

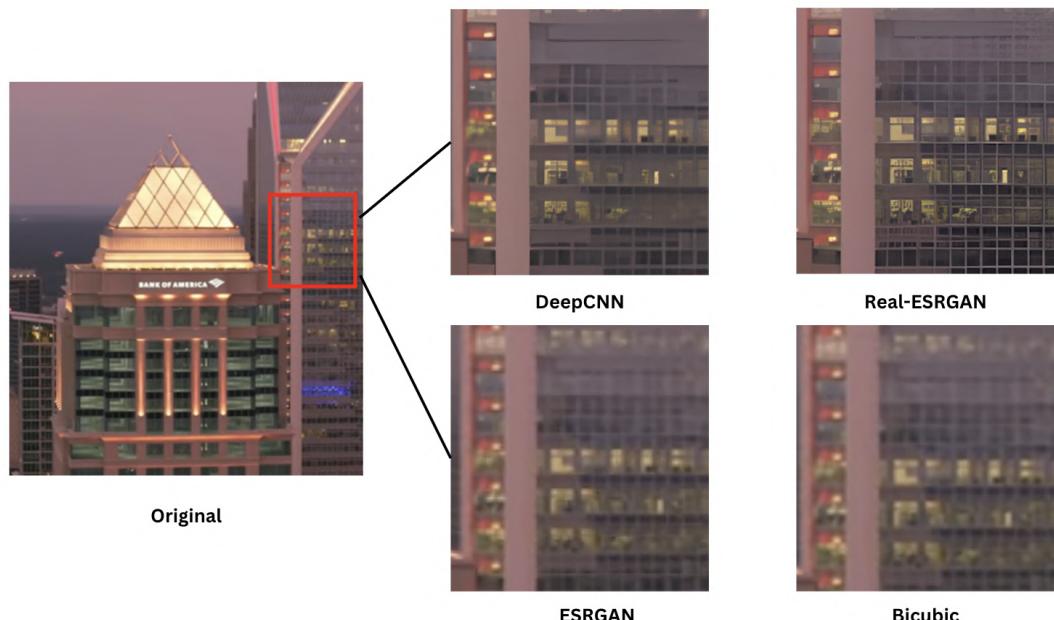


Figure 3.5: Analysis of models on building structure

3.2 Quantitative Analysis

3.2.1 Peak Signal To Noise Ratio

The values from Figure 3.1 tell us that bicubic interpolation offers slightly better reconstruction quality in terms of matching the overall signal strength, with a PSNR of 16.86 dB, followed closely by Real ESRGAN (16.63 dB), deepCNN (16.22 dB), and ESRGAN (16.08 dB).

3.2.2 Structural Similarity Index

SSIM scales from -1 to 1, with a value of 1 signifying complete similarity. According to the SSIM metrics, Real ESRGAN shows the highest structural similarity to the original images with a value of 0.23, indicating superior preservation of structural integrity compared to other techniques. Following closely, Bicubic interpolation registers an SSIM of 0.215. Meanwhile, deepCNN and ESRGAN record lower SSIM values of 0.17 and 0.15, respectively, suggesting less effectiveness in maintaining structural fidelity.

3.2.3 Learned Perceptual Image Patch Similarity (LPIPS)

Lower values of LPIPS indicate higher perceptual similarity. According to Figure 3.1, deepCNN and ESRGAN have equal performance (0.16 and 0.15 respectively), suggesting a closer perceptual match to the original images than bicubic interpolation (0.19). Real ESRGAN performs the best in this metric (0.14), indicating its superior capability in generating images that are perceptually closer to the original ones from a human perspective.

Model	PSNR (dB)	SSIM	LPIPS
DeepCNN	16.22	0.17	0.16
Bicubic	16.86	0.215	0.19
ESRGAN	16.08	0.15	0.15
Real ESRGAN	16.63	0.23	0.14

Table 3.1: Comparison of Image Enhancement Models

3.3 Assessing New Models on Drone Images

As a continuation of my mid-semester work, this part of the study introduces an evaluation of the latest super resolution models that have recently emerged as potential competitors to Real-ESRGAN. This section aims to expand the research scope by focusing on how well these newer models preserve textural details and overall image realism in drone captured images. The goal is to explore whether these advanced models provide any significant improvements over earlier technologies particularly in their handling of complex textures which is crucial for drone image enhancement. [11]

SwinIR: In practical tests, SwinIR demonstrated exceptional detail preservation, especially in architectural elements and natural landscapes. The model's effectiveness in maintaining clarity and fine details supports its suitability for precision-required imaging.

BSRGAN and A-ESRGAN: Observations showed that while these models improved over standard methods, they were less consistent in handling complex textures and often introduced minor artifacts. The artifact presence in densely textured areas reflects the challenges present in the adversarial training processes used by these models.

Real-ESRGAN: Real-ESRGAN was observed to maintain superior performance across various tests, effectively balancing detail enhancement and artifact suppression. This balance resulted in more natural and realistic image enhancements, confirming

its advanced capability in preserving the integrity of original textures while enhancing resolution.

Due to Real-ESRGAN's impressive performance in both visual quality and quantitative assessments (see Table 3.2), Real-ESRGAN was selected for additional fine-tuning, specifically aimed at enhancing drone-captured images. This process is intended to improve the model's effectiveness on this specific dataset, thereby boosting its applicability in real-world scenarios.

While SwinIR also demonstrated commendable performance, particularly in maintaining high-quality details in specific image types, Real-ESRGAN provided a more straightforward approach to fine-tuning. Additionally, it exhibited superior performance across a broader range of image types, making it the more practical choice for our comprehensive needs.

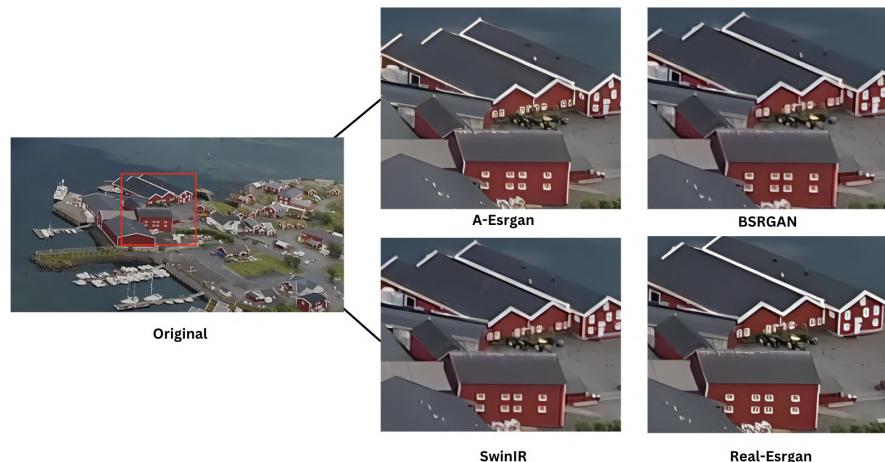


Figure 3.6: Analysis of models on window structures

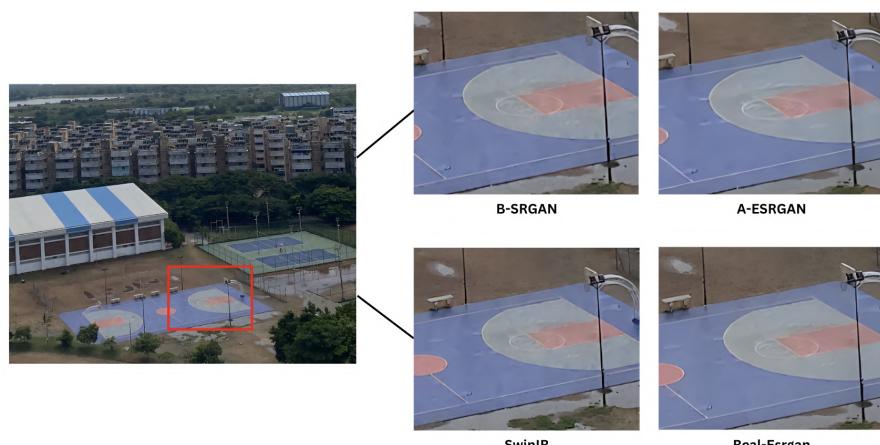


Figure 3.7: Analysis of models on flooring

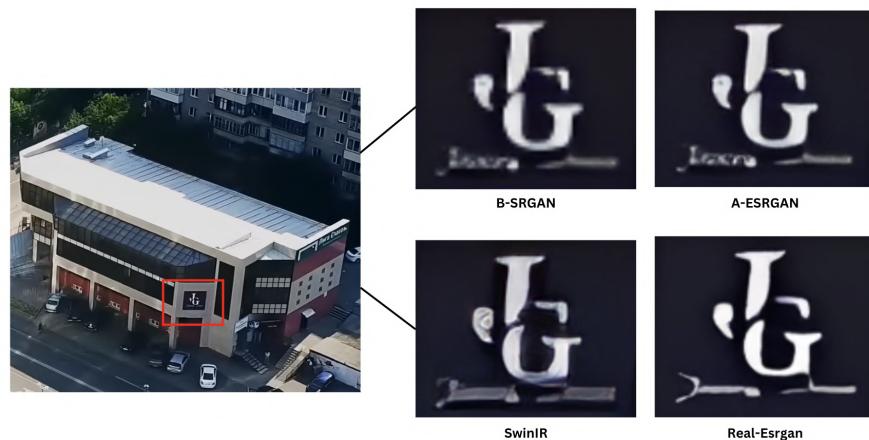


Figure 3.8: Analysis of models on letters

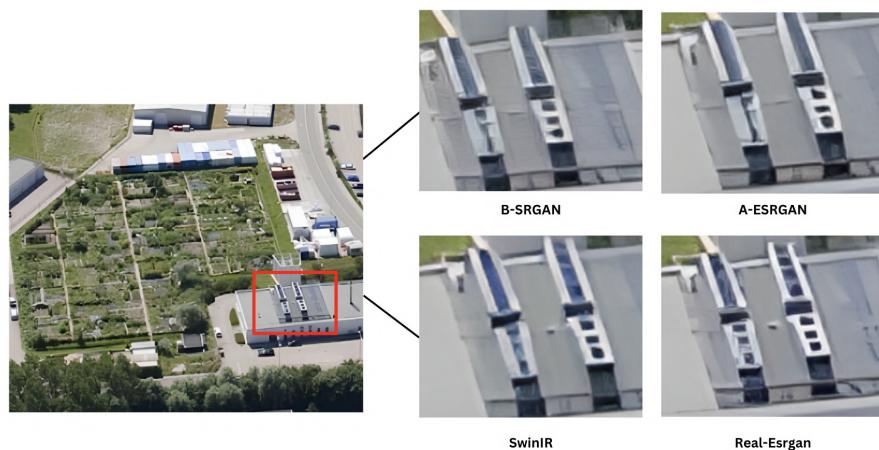


Figure 3.9: Analysis of models on fine details

Model	LPIPS	FSIM
Real ESRGAN	0.184	0.997
SwinIR	0.193	0.996
BSRGAN	0.422	0.942
A-ESRGAN	0.245	0.972

Table 3.2: Comparison of Super Resolution Models

3.4 Fine-Tuning Real-ESRGAN

3.4.1 Data Preparation

To begin the fine-tuning process, we gathered a diverse dataset of 700 images, crowd-sourced from various locations [12] [13], with a focus on buildings, trees, and landscapes. This dataset was tailored to represent the typical subjects encountered in drone photography. To prepare the images for training, we generated multi-scale versions by downsampling high-resolution (HR) images to create several ground-truth images at different scales. This approach helps in training the model to handle images of varying resolutions effectively. Additionally, the drone dataset images were cropped into smaller sub-images to speed up I/O and processing times during training.

3.4.2 Hyperparameter Adjustments

The learning rate was carefully reduced to allow for subtle adjustments to the pre-trained model, ensuring that the already learned features were not drastically altered. Initially, the model was trained for 10,000 iterations across 175 epochs to assess stability and performance. After confirming the model's stability, we extended the training to 100,000 iterations, allowing for a deeper and more comprehensive fine-tuning.

3.4.3 Augmentation Techniques

To better simulate real-world conditions, the images were subjected to various augmentation techniques. This included introducing scenarios with unclear daylight to test the model's performance under less than ideal lighting conditions. Different elevations were also considered to mimic the ones taken by drone. Color variations were also introduced to ensure the model could handle different lighting scenarios and color profiles effectively, making it more robust and versatile when applied to real-world drone imagery. (see Figure 3.10 , Figure 3.11)

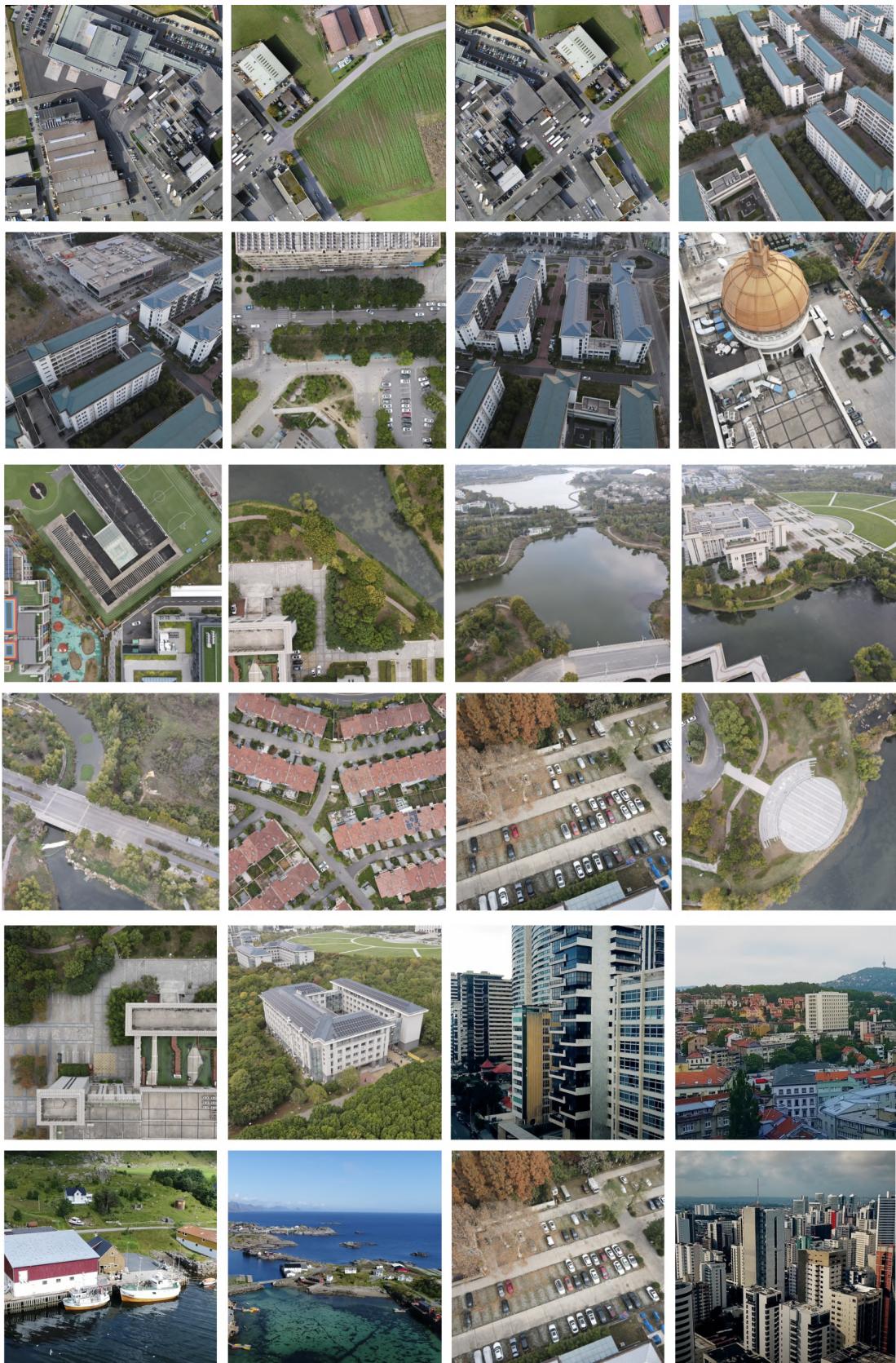


Figure 3.10: Dataset for finetuning

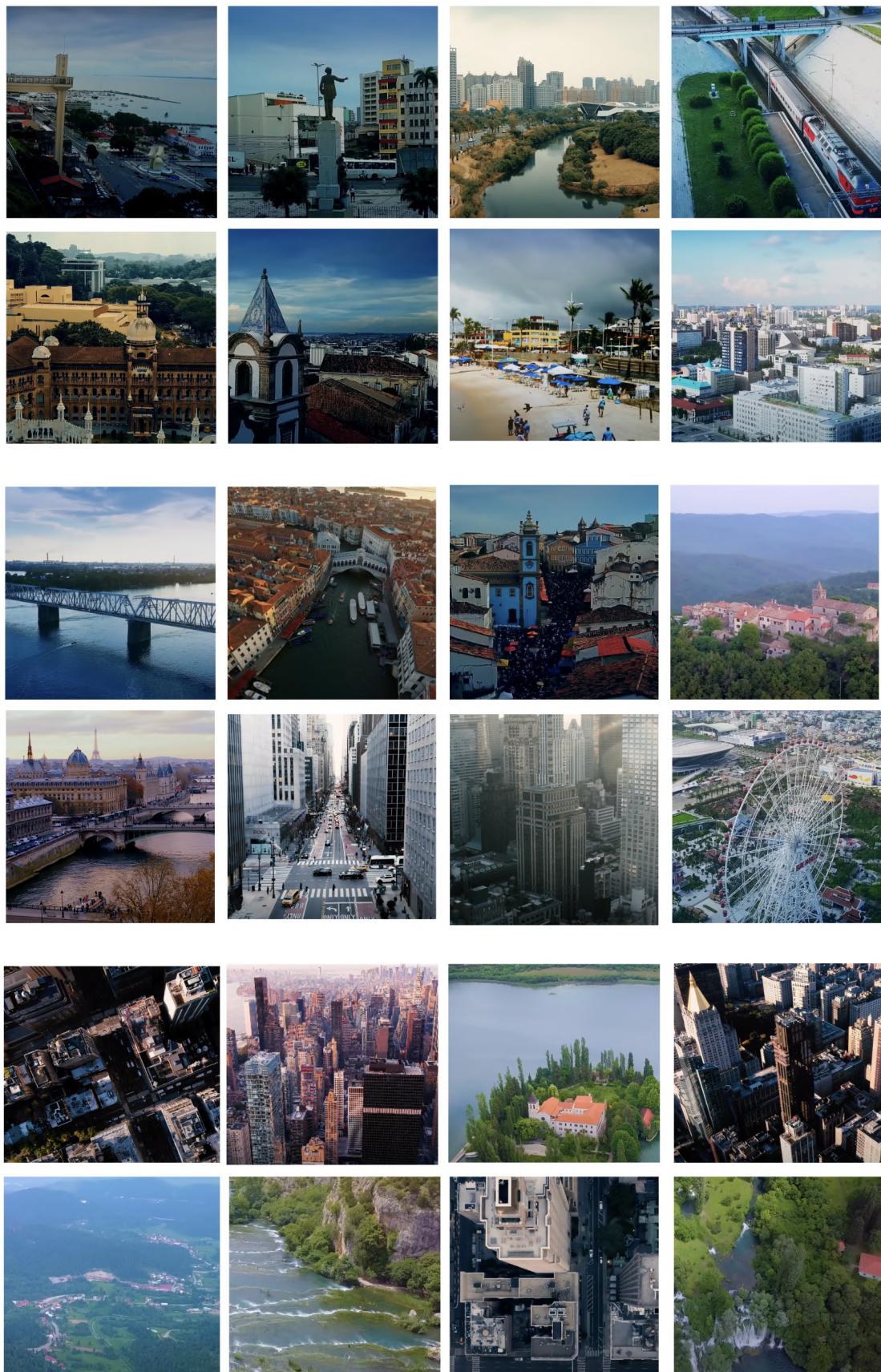


Figure 3.11: Finetuning dataset with variable lighting conditions and elevations

3.5 Outcomes of Fine-Tuning Real-ESRGAN

The process of fine-tuning Real-ESRGAN has enhanced the model's ability to capture and refine detailed textural and structural elements across diverse settings:

Architectural Edges: The fine-tuned model demonstrated improved definition in architectural features, particularly in the edges of walls and buildings. These features were rendered with greater smoothness and less distortion, contributing to a more precise and clearer representation of architectural forms. (see Figure 3.12)

Vegetation Textures: In natural landscapes such as fields, the grass texture appeared more realistic and detailed when processed by the refined model. This improvement is important for environmental monitoring applications where accurate texture reproduction can greatly influence the management of natural resources. (see Figure 3.13)

Handling Distorted Imagery: There was a notable enhancement in the model's ability to process images with high distortion, such as those depicting vehicles. The fine-tuned Real-ESRGAN was capable of delineating more distinct outlines and detailed features, essential for applications that demand high levels of detail for object identification and situational analysis. (see Figure 3.14)

Window Details in Buildings: The model's adjustments led to a significant enhancement in the precision of window details within building images. These improvements not only augment the visual aesthetics of the images but also enhance their utility for tasks that require detailed building inspections or real estate evaluations.(see Figure 3.16)

In addition to qualitative improvements, quantitative analysis was also conducted to evaluate the performance enhancements achieved by fine-tuning Real-ESRGAN. The analysis focused on two critical metrics: Learned Perceptual Image Patch Similarity and Feature Similarity Index. LPIPS measures the perceptual difference between the original and super-resolved images, offering insights into the model's ability to preserve image fidelity from a human visual perspective. FSIM, on the other hand, assesses the similarity in structural features and textural details, which are vital for applications requiring high levels of detail and accuracy. The results of this quantitative analysis

are presented in Figure 3.3 providing a comparative view of the original and fine-tuned models.

3.5.1 Remaining Challenges

Despite the significant enhancements observed, there are still areas requiring further development:

Subtle Textural Details: While the model considerably improved the rendering of major textural elements, the finer textures, particularly in low-light conditions or at extended distances, still lack adequate clarity and definition. Improving the model's capability to handle these finer details could expand its application range.

Complex Lighting Conditions: In scenarios featuring complex interplays of light and shadow, the model occasionally struggled to retain detailed information without introducing noise. Addressing this issue is essential for enhancing the model's performance in dynamically lit environments.



Figure 3.12: Analysis of models on building structure



Figure 3.13: Analysis of models on building structure

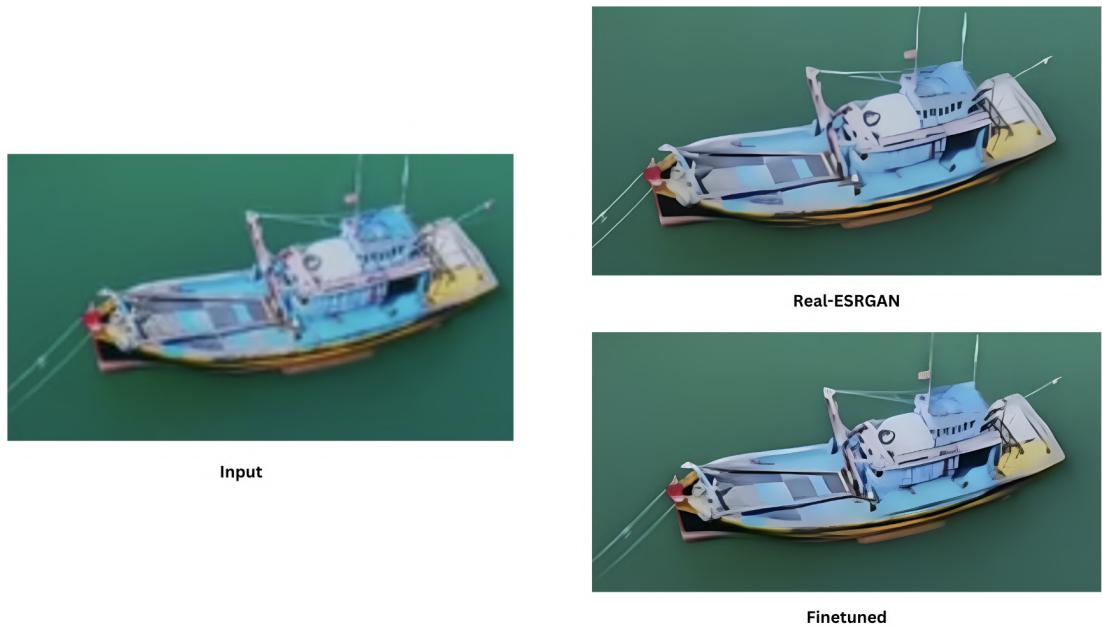


Figure 3.14: Analysis of models on distorted image



Figure 3.15: Analysis of models on distorted image



Figure 3.16: Analysis of models on building structure

Model	LPIPS	FSIM
Real ESRGAN	0.184	0.997
Finetuned Real ESRGAN	0.172	0.998

Table 3.3: Quantitative analysis of the two models

CHAPTER 4

CONCLUSION AND FUTURE SCOPE

4.1 Conclusion

This project has demonstrated noticeable improvements in image resolution by fine-tuning Real-ESRGAN, specifically for drone-captured images. The refinements in preserving details, reducing artifacts, and enhancing texture realism highlight the benefits of fine-tuning pre-trained models on targeted datasets. Particularly, Real-ESRGAN exhibited robust performance across both qualitative and quantitative evaluations, outperforming other models assessed in the study.

The fine-tuning process, which was adapted to specific environmental and architectural scenarios, notably enhanced the model's capability to process complex textures and lighting conditions. These enhancements are crucial for fields that demand high-quality image enhancement, such as urban planning, environmental monitoring, and sophisticated surveillance systems. This focused improvement approach ensures that Real-ESRGAN is well-suited to meet the needs of advanced applications requiring precise image details

4.2 Future Scope

4.2.1 Exploring New Models and Techniques

Future research can explore emerging and advanced models like Stable Diffusion, which offers state of the art capabilities in generating images through diffusion-based techniques. Integrating such models could provide alternative pathways to achieve even more realistic image enhancements with potentially lower computational costs.

4.2.2 Extended Applications

The application of advanced super-resolution models can be expanded to video enhancement, which remains largely unexplored. This could revolutionize areas like film restoration, live-event broadcasting, and real-time surveillance.

4.2.3 Enhancement of Model Training and Efficiency

There is substantial scope for improving the efficiency of these models to facilitate deployment on edge devices, which would be beneficial for real-time applications. Techniques such as model pruning, quantization, and the use of more efficient training algorithms could be investigated to enhance scalability and performance.

4.2.4 Further Fine-Tuning and Optimization

Continued fine-tuning of Real-ESRGAN and other models on specialized datasets with higher number of images can help in refining their capabilities to handle specific challenges such as high levels of noise and extreme variations in image quality.

By advancing these areas, future work can ensure that super-resolution technology not only continues to evolve but also becomes more integral and accessible for both industrial and consumer applications.

REFERENCES

- [1] R. Kellermann, T. Biehle, and L. Fischer, “Drones for parcel and passenger transportation: A literature review,” *Transportation Research Interdisciplinary Perspectives*, vol. 4, p. 100088, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2590198219300879>
- [2] T. Banu, G. Borlea, and B. Constantin, “The use of drones in forestry,” *Journal of Environmental Science and Engineering B*, vol. 5, 11 2016.
- [3] B. Siddappaji and K. Akhilesh, *Role of Cyber Security in Drone Technology*, 01 2020, pp. 169–178.
- [4] W. Rawat and Z. Wang, “Deep convolutional neural networks for image classification: A comprehensive review,” *Neural Computation*, vol. 29, no. 9, pp. 2352–2449, 2017.
- [5] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial networks,” 2014.
- [6] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. Change Loy, “Esrgan: Enhanced super-resolution generative adversarial networks,” in *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, September 2018.
- [7] X. Wang, L. Xie, C. Dong, and Y. Shan, “Real-esrgan: Training real-world blind super-resolution with pure synthetic data,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, October 2021, pp. 1905–1914.
- [8] Z. Wei, Y. Huang, Y. Chen, C. Zheng, and J. Gao, “A-esrgan: Training real-world blind super-resolution with attention u-net discriminators,” 2021.
- [9] J. Liang, J. Cao, G. Sun, K. Zhang, L. V. Gool, and R. Timofte, “Swinir: Image restoration using swin transformer,” 2021.
- [10] Kaggle and Contributors, “Stanford drone dataset,” <https://www.kaggle.com/datasets/aryashah2k/stanford-drone-dataset>.
- [11] X. Lin, B. Ozaydin, V. Vudit, M. E. Helou, and S. Süsstrunk, “Dsr: Towards drone image super-resolution,” 2022.
- [12] Image and E. Visual Representation Lab (IVRL), “Dsr dataset,” <https://github.com/IVRL/DSR>.
- [13] R. Robin, “Vdd: Visual design dataset,” <https://github.com/RussRobin/VDD>.