# A Deep Learning Framework for Semantic Colorization of Synthetic Aperture Radar Images

S.Varun[1], Sadiya[1], R.Akash[1], P.Praveen[1], K.Venkata Ramana[2], Dr. S Shiva Prasad[3*]

[1]Student, Department of CSE(DataScience), MallaReddy Engineering college, Secunderabad

[2] Assistant Professor, Department of CSE (Data Science), MallaReddy Engineering college, Secunderabad

[3*] Professor, Department of CSE (Data Science), MallaReddy Engineering college, Secunderabad

**Abstract- Synthetic Aperture Radar (SAR) image is useful in remote sensing because it is able to capture the images, regardless of weather and illumination conditions. But, SAR images are greyscale and naturally difficult to interpret by human or machine. In this paper, we propose a deep learning-based method for colorizing SAR images based on gray SAR and visual image pairs, with the goal of generating natural-looking color images from gray-level SAR images.**
**The solution proposed adopts a hybrid design, which leverages a Swin Transformer as the encoder for multi-scale feature extraction and utilizes an HR Net-based architecture as the decoder for high-resolution color generation. The model is learning cross-modal feature relationships in a supervised way, because it's trained on pairs of SAR–optical images. To maintain the structural consistency and improve the visual reality, Mean Squared Error (MSE) and perceptual loss functions are utilized. Experimental results indicate that the method is capable of generating good quality colorized SAR images, preserving structures and texture. This work contributes to an increased accessibility of SAR data for applications.**

**Keywords: Synthetic Aperture Radar (SAR), SAR Image Colorization, Deep Learning, Swin Transformer, HRNet, Cross-Modal Learning, SAR–Optical Image Pairing, Feature Extraction, Perceptual Loss, Mean Squared Error (MSE), Remote Sensing, Image-to-Image Translation.**

## I.INTRODUCTION

Synthetic Aperture Radar is almost the only means of remote sensing independent both of the sunlight and atmospheric conditions. This is quite different from passive optical sensors, as SAR systems emit microwave signals and receive reflected echoes back so that image data can be obtained even under cloud overcast, rain, or nighttime conditions. All these advantages make the application of SAR imaging particularly suitable for nearly every purpose: from environmental studies to disaster relief, from surveillance activities to terrain analysis.

However, weighed against these positives, the major drawback of SAR imagery is that it is captured in grayscale form. Lack of color information severely limits intuitive insight, making the interpretation by non-experts difficult. Among others, vegetation, water bodies, urban areas, and roads can hardly be visually distinguished. This is a serious impediment to the effective usage of SAR imagery in decision-making scenarios.

Although SAR imaging offers reliable and continuous data acquisition, the interpretation of SAR imagery remains problematic because of the interactive complexity of microwave signals with surface materials. Surface roughness, moisture content, geometry, and dielectric properties affect the back-scattered intensity in SAR images. Visual patterns are often punitive in comparison with optical imagery. Consequently, the absence of chromatic cues limits the direct usability of SAR data to analysts who are more familiar with the optical representations. Colorization of SAR images aims at bridging this semantic gap by translating the radar intensity patterns into meaningful color representations. Unlike conventional enhancement techniques, deep learning–based colorization models learn implicit correlations between grayscale textures and their corresponding color distributions from data. This data-driven approach allows for colorized outputs that enhance the visual clarity without losing the critical structural and textural information inherent in SAR imagery. Recent breakthroughs in attention-based architectures have significantly enhanced the modeling of long-range spatial dependencies in images. Transformer-based models, particularly, Swin Transformers, capture both the local and global contextual features

efficiently by hierarchical window attention mechanisms. This has increased the same benefits in SAR imagery, where terrain patterns and land-cover characteristics normally range in multiple spatial scales. On the other hand, high-resolution spatial detail must be preserved to save object boundary and fine textures, which are important for correct interpretation. In light of these requirements, the proposed framework incorporates a Swin Transformer encoder with a High-Resolution Network (Hr-net) decoder. This hybrid architecture enables effective multi-scale feature extraction while preserving high-resolution representations at each step in decoding. Coupling global contextual awareness with fine-grained spatial fidelity, the model generates the colorized SAR images with enhanced semantic consistency and visual realism.

Recent breakthroughs in deep learning, especially in image-to-image translation, have opened new horizons beyond such limitations. CNN and Transformer-based architectures have shown outstanding performance in learning complex spatial and contextual relationships between different domains of an image. This motivated us to propose a deep learning framework that leverages the strengths of Swin Transformers and HR-Net for performing colorization of SAR images. In this work, the ultimate goal is to provide photo-realistic and structurally consistent colorized SAR images that improve interpretability without loss of SAR-specific information. The rest of this work presents the related works, followed by the proposed methodology, experimental results, and conclusions.

## II.LITERATURE SURVEY

### LITERATURE SURVEY

Accordingly, image colorization has been a widely studied problem in the computer vision arena, especially under the grayscale to RGB image translation paradigm. The early approaches relied on heuristic-based pseudo-coloring or manual color mapping and yielded less visually consistent and semantically incorrect outputs. With the advent of deep learning, data-driven colorization of images that learn complex mappings between intensity patterns and color distributions was enabled.

Conditional generative adversarial networks greatly advanced image-to-image translation tasks. The Pix2Pix framework proposes effective paired image translation by jointly optimizing adversarial and reconstruction losses, making it one of the standard baselines for colorization applications. In order to reduce the dependency on paired datasets, Cycle GAN proposes cycle-consistency constraints that enable unpaired image translation and allow cross-domain learning.

While in the domain of remote sensing, deep learning techniques have been successfully conducted on Synthetic Aperture Radar image classification, segmentation, and de speckling, the colorization of SAR images remains relatively under explored. Due to inherent speckle noise, geometric distortions, and insufficient aligned SAR- optical datasets, traditional approaches to SAR colorization are often based on pseudo-coloring or CNN-based methods, which frequently produce low-resolution results with a lack of semantic realism.

Schmitt and Zhu demonstrated the difficulties of SAR–optical data fusion using the SAR optical data set by focusing on the modality gaps caused by the different imaging physics. Although unsupervised frameworks, such as Cycle GAN, allow training without paired data, supervised models offer in general much better fidelity results if aligned ground truth is available.

Recent development along various Transformer-based architectures has seen further improvement in the task of image-to-image translation. Swin Transformers capture global context with reduced computational complexity, while High-Resolution Networks (HR Net) preserve fine-grained spatial representations throughout the network. Although deep colorization networks proposed by Zhang et al. have demonstrated strong semantic awareness for natural images, they have not been optimized directly for SAR-specific noise characteristics. This work integrates Swin Transformers and HR Net in an effort to overcome these limitations to produce superior semantic consistency with improved visual quality in the colorization of SAR images.

Despite the progress achieved by existing image-to-image translation models, transferring color information from optical imagery to SAR data remains a challenging task due to the fundamental differences in sensing mechanisms. Optical images capture reflected visible light and naturally encode color semantics, whereas SAR images represent microwave back scatter responses that are heavily influenced by surface roughness, orientation, and material properties. This modality gap often leads to ambiguous mappings when conventional

colorization models are applied directly to SAR imagery. Consequently, there is a growing need for architectures that can effectively model both global contextual relationships and local structural details while being robust to SAR specific noise patterns. Hybrid frameworks that combine attention-based models with high-resolution feature preservation have shown promise in addressing these challenges, motivating further exploration into Transformer-guided SAR colorization techniques.

### III.PROPOSED METHODOLOGY

The methodology for colorizing grayscale Synthetic Aperture Radar images using a deep learning framework. The design approach here is essentially a supervised learning strategy based on paired SAR and optical images to learn cross-modal representations. It includes dataset preparation, preprocessing, architectural designs, loss formulation, and also other implementation details. A hybrid Swin Transformer and High-Resolution Network architecture is implemented to make an effort in capturing multi-scale contextual information while preserving fine spatial details. By carefully designing the preprocessing steps and optimizing the training procedures, the framework achieves realistic and well-structured colorized SAR images that fit well with practical applications for remote sensing.

### 3.1 Dataset Preparation

Dataset preparation is essential for effective supervised SAR image colorization. Only well-aligned SAR–optical image pairs with accurate spatial correspondence were selected to ensure reliable cross-modal learning. Misaligned, noisy, or low-quality samples were removed to avoid incorrect training. The dataset includes diverse land-cover scenes such as urban regions, vegetation, water bodies, and mixed terrains, helping the model learn generalized color mappings. Finally, the curated data was split into training, validation, and testing sets to prevent data leakage and support unbiased performance evaluation.

### 3.2 Dataset Preprocessing

Preprocessing was performed to stabilize training and improve model convergence. SAR and optical images were normalized to bring pixel values into a common numerical range, reducing gradient instability. Since SAR images are single-channel, they were replicated into three channels to match the required input format of the deep learning encoder. Data augmentation techniques such as horizontal/vertical flipping and controlled rotation were applied to improve robustness and generalization. Alignment checks were also conducted to ensure SAR–optical pairs maintained strict spatial consistency for accurate learning.

### 3.3 System Architecture

It should be pointed out that the heart of the above system consists of a PyTorch implementation of the SwinHRNet.

Encoder (Swin Transformer): The Swin transformer encoder we use is named swin_base_patch4_window7_224. Contrary to traditional CNNs, Swin Transformer uses the shifted window approach to extract features at a local and global level, which is more crucial for the context description of the terrain, for example, distinguishing water surfaces from flat ground regions based on SAR images. The features from the encoder are expressed at multiple scales, represented by x1, x2, and x3.

Decoder HR Net-The design of the decoder follows that of the High-Resolution Network. In this way, the multi-scale feature vectors are processed concurrently using different branches and preserve the high-resolution feature vectors. As a consequence, the fine details of the structures are also preserved in the resulting colorful images.

The proposed Swin-HRNet model represents a compromise between contextualization at a global level and a high level of detail about spatial reconstruction. The Swin-Transformer encoder is able to hierarchically process the input image and facilitates very efficient modeling of long-range relationships with relatively small computational complexity. As a result, the encoder is able to capture multi-scale information and the rich detail of a SAR image.

To complement this architecture, the HR Net decoder is able to preserve the detail of high-resolution feature streams throughout the reconstruction process. Unlike traditional architectures, which rely on down sampling and up sampling, the HR Net does not suffer from loss of detail during down sampling and is able to communicate throughout the parallel resolution branches, allowing it to preserve boundaries of

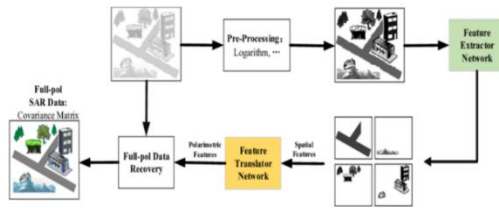objects, boundaries of terrain, and structural patterns throughout the reconstructed color images.



Figure 2. System Architecture

3.4 Loss Functions

To make the generated images mathematically accurate as well as perceptually realistic, the model minimizes a composite loss function:

MSE Loss: This measures the pixel-wise difference between the predicted color image and the ground truth optical image.

LPIPS - Perceptual Loss: We employ the Learned Perceptual Image Patch Similarity metric using a pre-trained VGG network. This loss ensures that the generated images match the ground truth in terms of high-level features and texture, other than just pixel intensity.

A hybrid pixel-level and perceptual loss function was thus adopted to guide the learning process toward results that are both visually coherent and structurally faithful. Although MSE loss enforces numerical similarity between predicted and ground-truth images, a single type of loss is not sufficient in capturing perceptual realism. Perceptual loss, based on deep feature representations, was therefore incorporated in this work, aimed at encouraging semantic content preservation and texture consistency. Jointly optimizing these components of loss allows the model to generate colorized output, both quantitatively and qualitatively, that aligns with human visual perception.
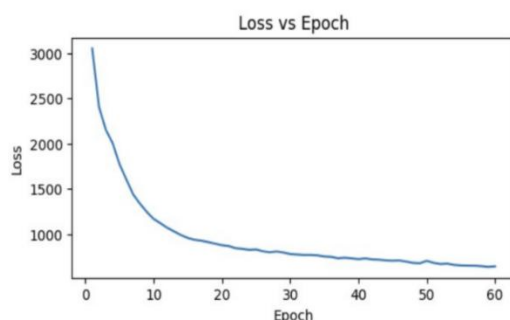


Figure 3 . Loss Function

3.5 Implementation Details

The model is trained using the Adam optimizer, with a learning rate of 1e-4. Training is done in an environment using CUDA acceleration, and this takes roughly 60 to 200 epochs16. To make the model deployable, a web application was implemented with Flask where users can upload a grayscale SAR image and download the colorized prediction through the browser interface.

Training was done on a GPU-accelerated environment in order to manage the high computational load required by Transformer-based architectures. Batch sizes were chosen according to the memory availability for stable training. Model checkpoints were saved after a certain period of time so that one could keep track of the performance and avoid overfitting. The model does the forward pass efficiently during inference, hence making the colorization almost real-time. The model has been integrated into a Flask-based web interface that demonstrates the model applicability interactively to present results of colorization of SAR images with interactive visualization without requiring technical proficiency.

IV..PROPOSED METHODS

Synthetic Aperture Radar (SAR) images contain rich structural details of terrains and objects but fail to exhibit natural colors in a meaningful manner for interpretation purposes. To overcome this drawback of SAR images, a new technology in this field is the application of deep learning algorithms for SAR image colorization. This project treats SAR image colorization as a supervised image-to-image translation problem.

The approach combines the latest developments in feature representation, leveraging Swin Transformers for capturing both local texture information as well as overall context, achieved via a shift-based window attention mechanism. Multiple scales of feature representation help to preserve small-scale structural information while being able to capture semantic context at different scales. High Resolution Network decoding encourages the maintenance of spatial information to ensure clean object boundaries and texture information within the colored output.

Cross-modal learning fills the gap between the SAR and Optical modalities to enable the network to predict semantically meaningful colors despite the

different sensing mechanisms. The training of the network is carried out using a hybrid loss function that ensures convergence by considering the Mean Squared Error and Perceptual loss.

These techniques combined provide a comprehensive framework which improves SAR colorization's fidelity, realism, and interpretation. By leveraging cutting-edge feature extraction via transformers, multiscales, and high resolutions, this proposed technique maintains both structural and semantic consistency. This proposed methodology offers a strong foundation for applications in remote sensing, environmental surveillance, and urban planning, where SAR interpretation's fidelity and realism are extremely important.

4.1.Supervised Image-to-Image Translation

Supervised Image-to-Image Translation: The task of coloring SAR images is made into a learning task where the grayscale SAR images and the optical image are paired, and the network is trained on these pairs of datasets. The need for paired data makes it possible for the network to learn the accurate mapping between the reflectivity pattern of the SAR image and the natural color information.

The advantage of supervised learning is that it enables end-to-end optimization, wherein the network itself tries to optimize the difference between the estimated and actual optical images. Thus, the colorization task becomes more robust than the case of unsupervised or rule-based techniques, which frequently face issues related to inconsistencies and unreasonable colorization artifacts. Additionally, the supervised translation learning mechanism is quite efficient and robust to the terrains of differences related to sensor noise levels, resolutions, and terrains.

Furthermore, this approach enables a quantitative analysis of the performance of the model based on the peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM). These parameters help in measuring the proximity of the correctly predicted color images with the actual optical images. Thus, supervised image-to-image translation is a well-grounded platform in the SAR colorization task containing elements like structural consistency, as well as visually realistic consistency.
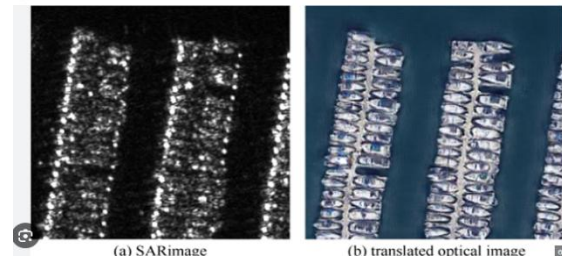


(a) SARimage                     (b) translated optical image

Figure 4. Image-to-Image Translation

4.2.Swin Transformer-Based Feature Extraction

The Swin Transformer is a powerful feature extractor, capable of learning local information and global context from SAR images simultaneously. The main strength of this approach is its shifted-window attention module, which enables the model to focus on different parts of the feature map without overlapping while maintaining efficiency. This is especially helpful in radar images, where the model needs to identify local textures as well as larger terrain patterns.

Contrary to traditional convolutional neural networks, the relationships between distant pixels are modeled dynamically by the Swin Transformers. This is very effective in the case of SAR images because the patterns formed by radar back scatter are spatially complicated. With the hierarchy of feature representations modeled by the transformer, the dependencies on various scales are utilized to emphasize the task of colorization.

Additionally, it is important to note that the Swin Transformer architecture is scalable when dealing with high-resolution SAR images. The capability of handling both high and low-resolution information promotes edge and object details preservation within the network. The incorporation of this model within a SAR colorization framework will ensure that features extracted are rich and meaningful enough to serve as a foundation for SAR image reconstruction.
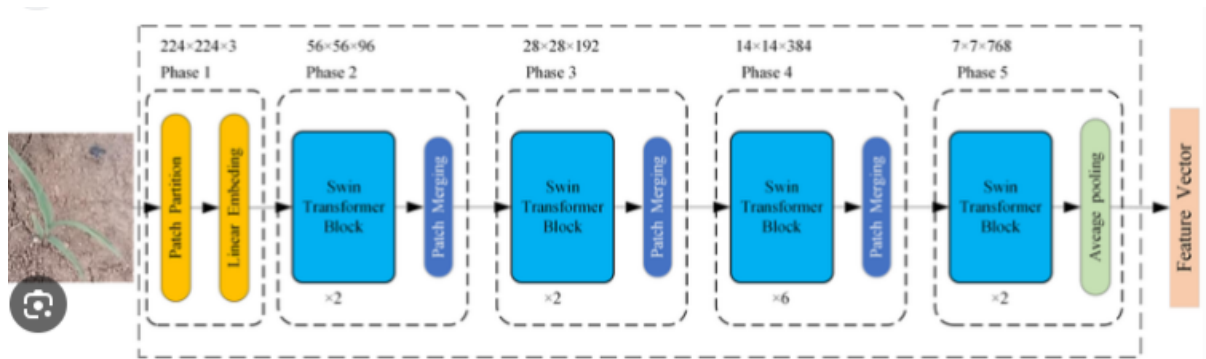
Figure 5. Swin Transformer-Based Feature Extraction

### 4.3. Multi-Scale Feature Representation

Multiple scales of feature representation play a vital role in capturing the varying spatial details that SR images contain. Using the network to produce different scales of maps enables it to capture local texture patterns and abstract semantic meanings simultaneously. Low-scale maps capture large-scale context dependencies and preserve small-scale details like edges and boundaries in the maps.

This helps the network in integrating the global and the local features of the image in the process of colorization. In the case of the SAR image, it is likely to have complex backscatters due to either the rugged nature of the surface or the presence of anthropogenic structures, and hence it demands multiresolution analysis to achieve realistic color transfer.

### 4.4. High-Resolution Network (HRNet) Decoding

In contrast to traditional decoders, which progressively down-sample and up-sample the feature maps, the decoder HRNet maintains parallel streams of high resolution, allowing the network to restore fine details without losing spatial resolution. This is critical for SAR images, where object boundaries and textural features are important.

Therefore, all the outputs would have both structural accuracy and semantic consistency because HR-Net maintains several parallel resolutions and fuses them by repeated exchange blocks. It will prevent blurring or color bleeding problems that happen usually in standard encoder-decoder architectures. It also enhances subtle terrain variations, increasing the model's representation capacity for applications in urban mapping or environmental monitoring.

Besides, the high-resolution decoding of HR-Net enhances feature extraction by incorporating multi-scale features comprehensively. While reconstructing the preserved spatial details, colorization would not affect geometric and structural integrity. So, the HR Net-based decoding generates visually coherent colorizations of SAR that are also faithful to the underlying scene.

### 4.5. Cross-Modal Learning Between SAR and Optical Data

Cross-modal learning fills the gap between SAR radar images and optical images, which are based on distinct principles of perception. SAR radar images perceive the roughness and dielectric constant of the material, whereas optical images perceive the reflected light within the visible spectrum. Cross-modal learning helps in determining the meaningful mapping for the model to efficiently apply colorization.

Through learning common latent features, it is possible for the network to deduce optical colors from radar backscatter values. Not only is this useful for realism, but it can improve semantic interpretation of the scene as well, given the fact that varied materials have a unique reflectivity value. Cross-modal learning prevents any inconsistencies in coloring, such as water, grass, and cities being colored incorrectly.
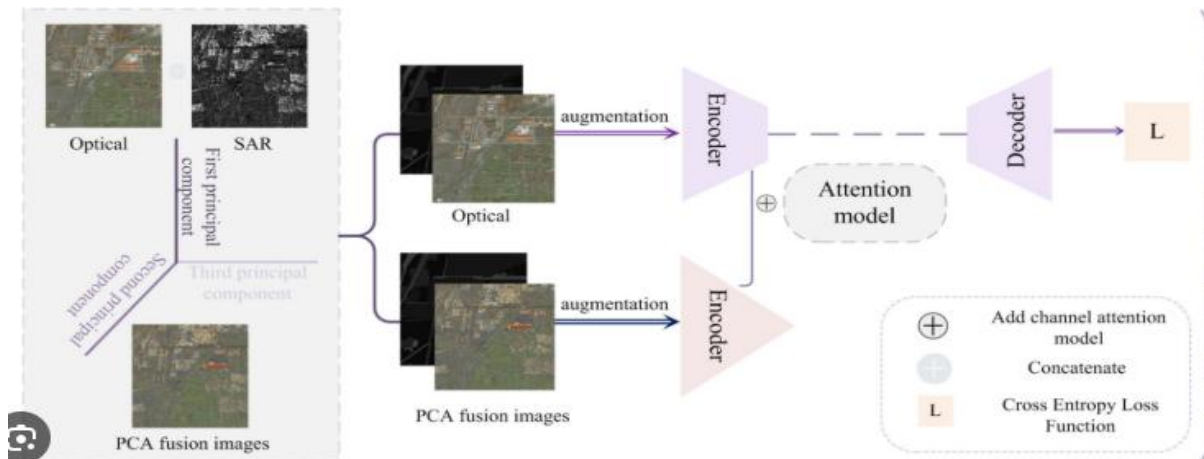
Figure 6. Cross-Modal Learning Between SAR and Optical Data

4.6.Composite Loss Optimization (MSE + Perceptual Loss)

To achieve composite loss, both pixel-wise loss and perceptual loss are used. Mean Squared Error (MSE) loss is used to ensure the predicted pixel values are close to the actual value of the image. Meanwhile, MSE sometimes produces smooth images with no realistic texture. To overcome this limitation, MSE loss can be combined with perceptual loss. Perceptual loss involves comparisons in terms of higher-level features derived from pre-trained models.

This joint loss function serves to combine low-level accuracy and high-level perception. The loss function consisting of MSE pursues high-precision reconstruction of the object boundary and details, and at the same time, promotes natural color transition and texture continuity. The particularity of this combination is more evident in SAR colorization: there is a natural ambiguity of correspondence between gray backscatters and color optical signals.

V..RESULTS

5.1. Training Performance
This proposed network was trained for a period of 60 epochs. In this case, it was evident that there was stable convergence of the optimization process. This is because the loss during training showed stable convergence without exhibiting any signs of divergence or oscillating. This is an indicative sign that there was a proper balance between the learning rate and the architecture used. This signified that it managed to learn the complex relationship between SAR and optical imagery.
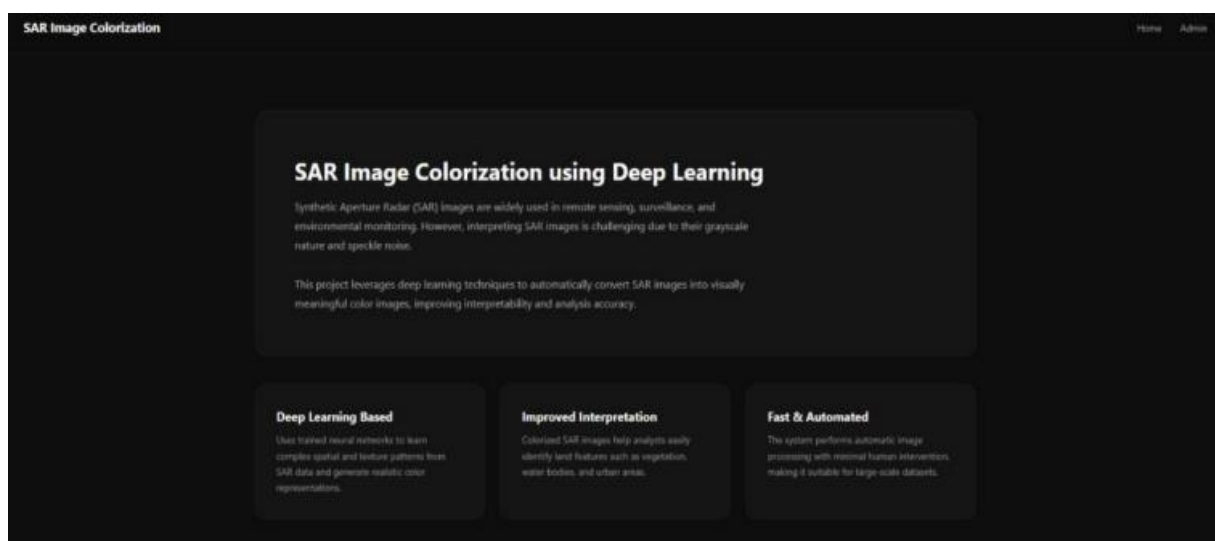


Figure 8. Interface of project

Besides, there are no symptoms of overfitting, such as erratic fluctuations of the loss or early stagnation, which further strengthens the robustness of the training. Such gradual reduction of the loss values reflects the capability of the model to generalize feature representations rather than memorizing the training samples. In general, the training dynamics confirm the effectiveness of the proposed approach for learning a meaningful cross-domain mapping between radar and optical image spaces.

5.2.Qualitative Analysis

For assessing the capability of generalization of the trained model, some qualitative experiments were performed on some untrained SAR images. These images used as inputs possessed some critical characteristics, which include speckle noise, low contrast, and intensity distributions of gray levels, and normally, images with these features possess difficulties in interpretation. Despite these difficulties, the model produced meaningful colorized images.

The different classes of land cover were well distinguished by distinct features in the prediction images. Vegetated regions were properly highlighted using natural shades of green, and water features were effectively shown using dark blue or black hues to easily distinguish them from the surrounding ground. Land and the built environment properly transitioned from one to another. These images used as inputs possessed some critical characteristics, which include speckle noise, low contrast, and intensity distributions of gray levels, and normally, images with these features possess difficulties in interpretation.

In contrast to the classical pseudo-coloring approach based on a fixed intensity-to-color mapping, in the proposed method the boundary of the object became sharper and fewer colors bled into each other. This shows that the model is not only able to apply a semantic color assignment but also to derive semantic information at a higher level from the SAR image.
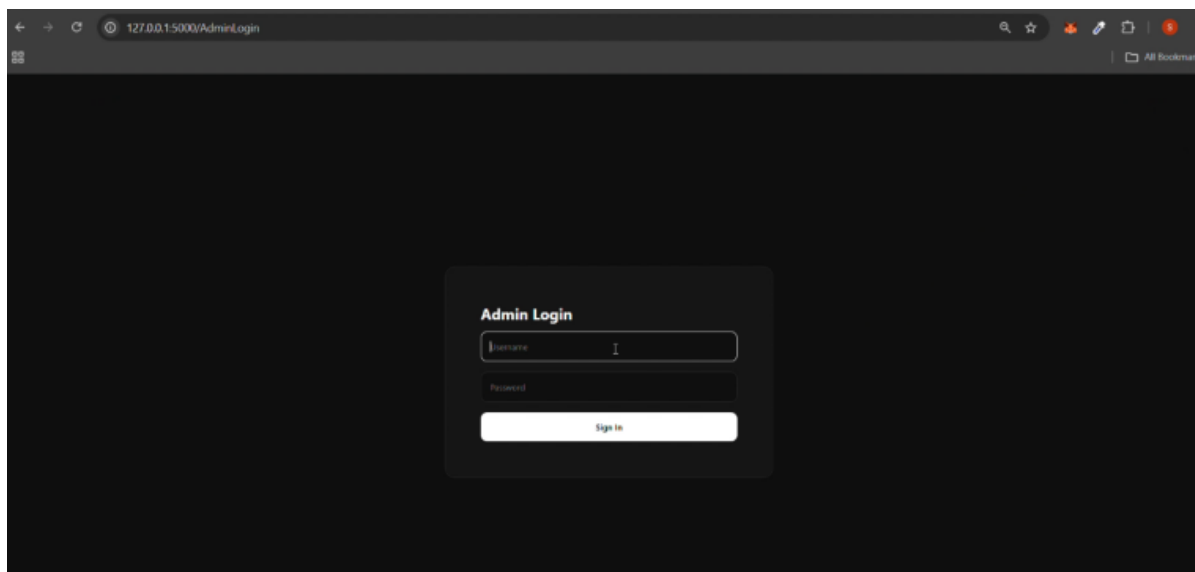


Figure 9. login page

5.3.Web Deployment For the ability of the proposed framework to be truly demonstrated, the resulting model is integrated with a Flask web interface, which makes it applicable and usable in real-world situations like image analysis and processing tasks. The interface is such that it is able to process and colorize the SAR image uploaded by the user instantaneously, and the colored image is stored under timestamped names, which ensures easier management of the resulting output. The web interface provides both the original and colored image, making it possible for the user to evaluate the effectiveness of the image processing task performed on the input image correctly, which makes the proposed framework applicable in image analysis and processing tasks like remote sensing and image processing.
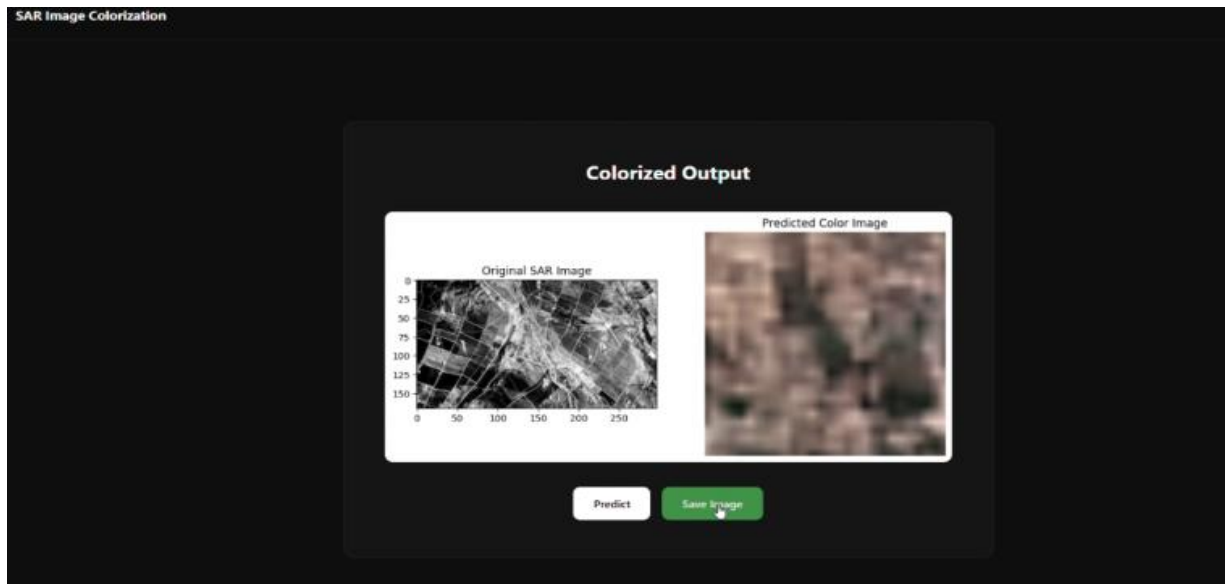
Figure 10. final output with download option

## VI.CONCLUSION

In this, the effectiveness of deep learning for the challenging task of Synthetic Aperture Radar (SAR) image colorization. A hybrid framework using a Swin Transformer encoder and an HRNet decoder is proposed to capture global contextual features while preserving fine spatial details. Unlike traditional pseudo-coloring methods, the proposed model is fully data-driven and adapts well to complex scene variations. Experimental results show that the system generates visually realistic and structurally consistent colorized SAR images even under speckle noise. The improved interpretability makes SAR imagery easier to understand for non-experts while retaining useful information for specialists. This enhanced visualization can support applications such as agriculture monitoring, land-cover analysis, environmental change detection, disaster assessment, and security surveillance. However, the approach relies on paired SAR–optical datasets, which limits scalability across regions and sensors. Future work will focus on unsupervised or weakly supervised learning methods such as CycleGAN and contrastive learning to use unpaired data. Additionally, real-time and edge deployment optimization through model compression and lightweight attention is planned. Overall, the framework offers a robust and scalable solution for making SAR data more accessible and interpretable in real-world remote sensing systems.

## REFERENCES

[1] Zhang, R., Isola, P., & Efros, A. A. (2016). *Colorful Image Colorization*. In Proceedings of the European Conference on Computer Vision (ECCV), 649–666.

[2] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... & Bengio, Y. (2014). *Generative Adversarial Nets*. Advances in Neural Information Processing Systems (NeurIPS), 2672–2680.

[3] Schmitt, M., Zhu, X. X. (2016). *SAR-Optical Data Fusion: State-of-the-Art and Future Challenges*. IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, 9(10), 4598–4623.

[4] Ronneberger, O., Fischer, P., & Brox, T. (2015). *U-Net: Convolutional Networks for Biomedical Image Segmentation*. In International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI), 234–241.

[5] Zhu, J. Y., Park, T., Isola, P., & Efros, A. A. (2017). *Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks*. IEEE International Conference on Computer Vision (ICCV), 2242–2251.

[6] Sentinel-1 and Sentinel-2 Data Products. European Space Agency (ESA). Retrieved from https://sentinel.esa.int

[7] Lee, J. S. (1981). *Refined Filtering of Image Noise Using Local Statistics*. Computer Graphics and Image Processing, 15(4), 380–389.

[8] Wang, Z., Bovik, A. C., Sheikh, H. R., & Simoncelli, E. P. (2004). *Image Quality Assessment: From Error Visibility to Structural Similarity*. IEEE Transactions on Image Processing, 13(4), 600–612.

[9] (2006). Anomaly Based Network Intrusion Detection with Unsupervised Outlier Detection. *IEEE International Conference on Communications*, 2006, 2388–2393. https://doi.org/10.1109/ICC.2006.254888