# TRANSFORMER-BASED SAR IMAGE DESPECKLING

*Malsha V. Perera, Wele Gedara Chaminda Bandara, Jeya Maria Jose Valanarasu, and Vishal M. Patel*

Johns Hopkins University
Department of Electrical and Computer Engineering
{jperera4, wbandar1, jvalana1, vpatel36}@jhu.edu

## ABSTRACT

Synthetic Aperture Radar (SAR) images are usually degraded by a multiplicative noise known as speckle which makes processing and interpretation of SAR images difficult. In this paper, we introduce a transformer-based network for SAR image despeckling. The proposed despeckling network comprises of a transformer-based encoder which allows the network to learn global dependencies between different image regions - aiding in better despeckling. The network is trained end-to-end with synthetically generated speckled images using a composite loss function. Experiments show that the proposed method achieves significant improvements over traditional and convolutional neural network-based despeckling methods on both synthetic and real SAR images. Our code is available at : `https://github.com/malshaV/sar_transformer`

***Index Terms***— Synthetic Aperture Radar, transformers, speckle, denoising

## 1. INTRODUCTION

Synthetic Aperture Radar (SAR), like other coherent imaging systems, is also affected by speckle which is a signal-dependent, spatially correlated noise. The presence of speckle impairs downstream tasks of SAR images such as segmentation, recognition, and detection. Therefore, despeckling of SAR images positively impacts these downstream tasks.

In the past few decades, different despeckling approaches have been proposed in the literature. Filter-based approaches used for despeckling can be generally categorized as local and non-local filters. Lee filter [1] and Kuan filter [2] are some despeckling methods which use local filters, while PPB [3] and SAR-BM3D [4] are examples for non-local filter-based despeckling methods. A survey of different SAR image despecking methods can be found in [5].

With the advent of deep learning, significant advances have been made in SAR despeckling. Unlike the despeckling approaches mentioned earlier, convolutional neural networks (CNNs) require a pair of noisy and clean (ground truth) images so that they can be trained in a supervised way. However, clean references for SAR images do not exist. Therefore, supervised CNN-based methods generate reference images via synthetic speckle generation or multi-temporal fusion. SAR-CNN [6] is a CNN-based despeckling network which is trained on SAR images transformed to the hormomophic form. The ground truth for SAR-CNN is generated via multi-temporal fusion. Wang *et al.* [7] proposed ID-CNN which directly estimates the noise in the original domain and is trained on synthetic SAR images. Here, the despeckled image is obtained by dividing the speckled SAR image by the estimated noise. Despeckling approaches such as [8, 9] proposed slight variations on the above CNN-based methods by introducing different architectures and loss functions.

With the success of transformers in natural language processing, they have been successfully adopted to many vision tasks in recent years. Dosovitskiy *et al.* [10] proposed Vision Transformer (ViT) which resulted in an impressive image classification performance on ImageNet as transformers have the ability to learn long-range dependencies in an image. In ViT, the input image is split into multiple linearly embedded patches which are fed into a transformer encoder. Following ViT, many studies have employed transformers in various computer vision tasks achieving state-of-the-art performance. For example, U-former [11] and SwinIR [12] are transformer-based methods proposed for image restoration. To the best of our knowledge, transformer-based SAR despeckling has not been studied in the literature.

To this end, we propose a transformer-based network for SAR image despeckling. The proposed network consists of a transformer encoder and a CNN decoder. Unlike ViT which only produces feature maps with fixed resolutions, the hierarchical nature of our proposed transformer encoder allows us to generate high resolution fine features as well as low resolution coarse features required for SAR image despeckling. Similar to [13], our transformer encoder avoids interpolating positional codes when performing inference. Therefore, the encoder can easily adapt to different test resolutions without

impairing the performance. The proposed network is trained end-to-end with synthetically speckled optical images. Finally, we compare the performance of our proposed network on synthetic images as well as real SAR images with several non-local filter and CNN-based methods where we achieve state-of-the-art performance.

## 2. PROPOSED METHOD

### 2.1. Speckle in SAR

For a SAR image with an average number of $L$ looks, the observed SAR intensity $y$ is related to the speckle-free SAR intensity $x$, as follows:

$$y = xn, \tag{1}$$

where $n$ is the multiplicative speckle. Under the hypothesis of fully developed speckle, $n$ follows a Gamma distribution with unit mean and a variance of $1/L$. Therefore, the probability distribution of $n$ is given by,

$$p(n) = \frac{1}{\Gamma(n)} L^L n^{L-1} e^{-Ln}, \tag{2}$$

where $\Gamma(.)$ is the Gamma function. Given $y$, our goal is to estimate $x$.

### 2.2. Network architecture

The proposed network architecture is illustrated in Fig. 1. The noisy input image $y$ is passed through an encoder comprising of 5 stages, each stage containing an overlap patch embedding block and a transformer block. The last stage of the encoder is followed by a convolutional projection block that acts as a decoder. At each encoder stage, the resolution of the input is halved to allow the transformer to learn both coarse and fine details. Output from each encoder stage (except for the final stage) is then fed to both the next encoder stage and the convolutional projection block. The output of the final encoder stage is passed only to the convolutional projection block. The components of the network are described in detail below.

**Overlap Patch Embedding Block:** Input to each encoder stage is first passed through an overlap patch embedding (OPE) block which uses the overlap patch merging process introduced in [13]. Its purpose is to combine overlapping feature patches to obtain features of the same size as that of non-overlapping patches before passing the features to the transformer block. In order to obtain the overlapped feature patches, the input to OPE block is passed through a convolutional layer of kernel size $k$, embedded dimensions (i.e. number of filters) $e$, stride $s$, and padding $p = k/2$. In this work, we set $s = 2$ for all the OPE blocks in the network. The $e$ values were set to 32, 64, 128, 320, 512 and $k$ values were set equal to 7, 3, 3, 3, 3 in each OPE block of stages 1 to 5, respectively. Next, the flattened output of the convolutional layer is followed by layer normalization.

**Transformer Block.** The output of each transformer block $T(\mathbf{I}_{in})$ with respect to an input $\mathbf{I}_{in}$ can be summarized as follows:

$$T(\mathbf{I}_{in}) = \text{MLP}(\text{DWC}(X(\mathbf{I}_{in}))) + X(\mathbf{I}_{in}) \tag{3}$$

where,

$$X(\mathbf{I}_{in}) = \text{MHA}(\mathbf{I}_{in}) + \mathbf{I}_{in}, \tag{4}$$

where, MHA, DWC and MLP correspond to multi-head attention layer, depth-wise convolution and multi-layer perceptron, respectively. Note that we perform layer normalization on $\mathbf{I}_{in}$ and $X(\mathbf{I}_{in})$ before passing those to MHA and DWC, respectively. In the original multi-head self-attention process, the queries $\mathbf{Q}$, keys $\mathbf{K}$ and values $\mathbf{V}$ have the same dimensions $d$, and the self-attention is calculated as:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right)\mathbf{V} \tag{5}$$

In the proposed network, each transformer block in stages 1 to 5 has 1, 1, 2, 4, and 8 attention heads, respectively. To reduce the computational complexity, we use the reduction ratio $R$ introduced in [14], and we set $R = 2$ for all the transformer blocks in the network. The self-attention features are then passed through the depth-wise convolutional block that provides positional information for transformers [13]. Subsequently, the output from the depth-wise convolution block is sent through a Gaussian error linear unit (GELU) before passing it through an MLP that comprises of a dropout layer and a linear layer. The output size of the MLPs in stages 1 to 5 are set to 32, 64, 128, 320, and 512, respectively.

**Convolutional Projection Block.** The convolutional projection block is used to upsample the outputs from the transformer blocks to the original image size as illustrated in Fig. 1. The upsampling layers increase the resolution by a factor of two. The output of the residual block (RB) for a given input $\mathbf{I}_{in}$ can be computed as follows:

$$\text{RB}(\mathbf{I}_{in}) = \text{Conv}_{3\times3}(\text{ReLU}(\text{Conv}_{3\times3}(\mathbf{I}_{in}))) + \mathbf{I}_{in}, \tag{6}$$

where $\text{Conv}_{3\times3}$ refers to $3 \times 3$ convolution layer and ReLU denotes rectified linear unit.

### 2.3. Loss Function

We train the proposed network using the $l_2$ loss ($\mathcal{L}_{l_2}$) given by:

$$\mathcal{L}_{l_2} = \|\hat{x} - x\|_2^2, \tag{7}$$

where $x$ and $\hat{x}$ denote the ground truth and the predicted output, respectively. In addition, we employ total variation loss in order to encourage smoothness while preserving edges. The total variation loss ($\mathcal{L}_{tv}$) is defined as follows

$$\mathcal{L}_{tv} = \sum_{i,j} |\hat{x}_{i+1,j} - \hat{x}_{i,j}| + |\hat{x}_{i,j+1} - \hat{x}_{i,j}|. \tag{8}$$

The overall loss function $\mathcal{L}$ is given by:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{l_2} + \lambda_2 \mathcal{L}_{tv}, \tag{9}$$

where $\lambda_1, \lambda_2$ are the weights that specify the contribution of $\mathcal{L}_{l_2}$ and $\mathcal{L}_{tv}$, respectively. In this study, we set $\lambda_1 = 1$ and $\lambda_2 = 5 \times 10^{-5}$ to put a strong emphasis on the $\mathcal{L}_{l_2}$ loss.
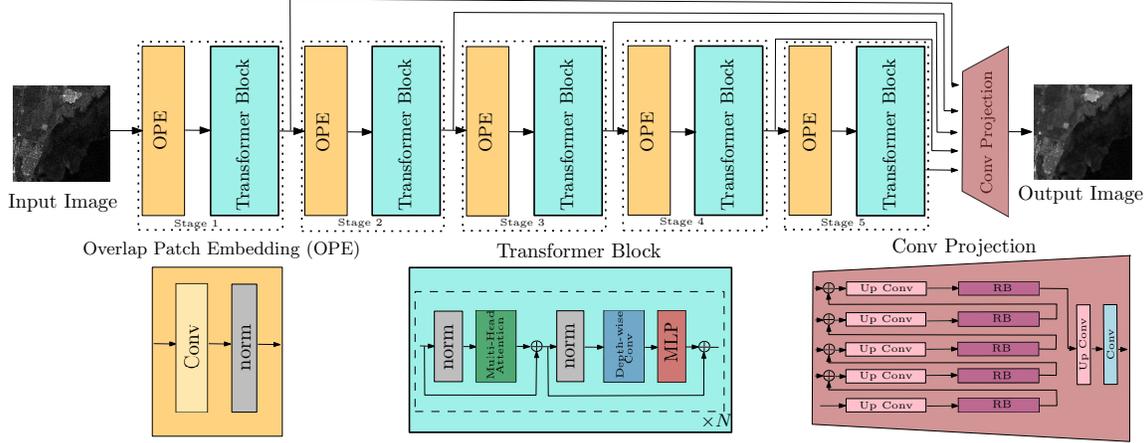
**Fig. 1**. Overview of the proposed despeckling transformer.

**Table 1**. Results on synthetic images of Set12 dataset [16].

| Method | PSNR | SSIM |
|---|---|---|
| PPB [3] | 21.90 | 0.599 |
| SAR-BM3D [4] | 23.51 | 0.701 |
| SAR-CNN [6] | 24.51 | 0.651 |
| ID-CNN [7] | 24.44 | 0.685 |
| Proposed method | **24.56** | **0.718** |

**Table 2**. Results on real SAR images.

| Method | Region 1 | | Region 2 | | Region 3 | | Region 4 | |
|---|---|---|---|---|---|---|---|---|
| | ENL ↑ | Cx ↓ | ENL↑ | Cx↓ | ENL↑ | Cx↓ | ENL↑ | Cx↓ |
| PPB | 87.0 | 0.11 | 125.5 | 0.09 | 21.0 | 0.22 | 117.8 | 0.09 |
| SARBM3D | 110.8 | 0.09 | 104.2 | 0.10 | 34.9 | 0.17 | 122.9 | 0.09 |
| SARCNN | 87.4 | 0.11 | 51.12 | 0.14 | 30.7 | 0.18 | 68.6 | 0.12 |
| IDCNN | 47.6 | 0.14 | 33.3 | 0.17 | 23.1 | 0.21 | 34.5 | 0.17 |
| Proposed | **154.4** | **0.08** | **171.6** | **0.08** | **39.2** | **0.16** | **133.19** | **0.08** |

## 3. EXPERIMENTS AND RESULTS

To train the proposed network, we generated single-look ($L = 1$) synthetic speckle images using optical images following equations 1 and 2. We used the Berkeley segmentation dataset (BSD) [15] to obtain the optical images where 450 and 50 images of size $256 \times 256$ were allocated for training and validation, respectively. The proposed network was implemented using PyTorch and trained with a learning rate of 0.0002 for 400 epochs. The performance of the proposed algorithm on synthetic speckled images were tested using a set of well-known testing images [16] in terms of Peak Signal-to-Noise ratio (PSNR) and Structured Similarity Index (SSIM). We have compared the performance of the proposed network with PPB [3], SAR-BM3D [4], SAR-CNN [6] and ID-CNN [7]. Note that the two CNN-based methods (SAR-CNN and ID-CNN) were trained on the same synthetic data as the proposed method.

The corresponding results are summarized in Table 1. From Table 1, it can be observed that the proposed method outperforms both traditional and CNN-based despeckling methods in terms of PSNR and SSIM when tested on synthetically speckled images.

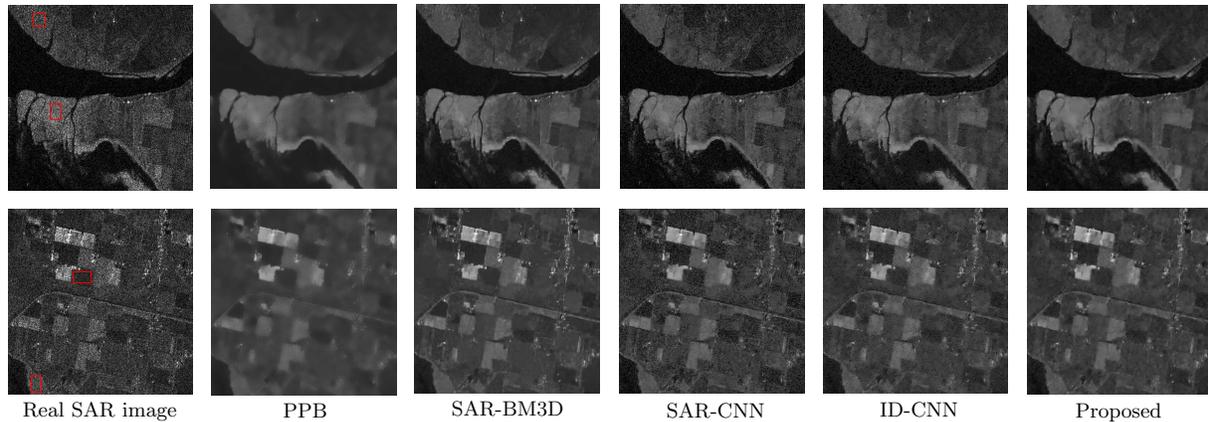We also evaluate the despeckling performance of our proposed method by testing on real SAR images. We compare the performance with the same despeckling approaches explained above. The despeckled results on the real SAR images are visualized in Fig. 2 for qualitative comparison. Since real SAR images do not have clean ground truth images, we use equivalent number of looks (ENL) and the coefficient of variation (Cx) derived in a homogeneous region (regions considered in this study are marked as red boxes in Fig. 2) as quantitative measures for comparison. ENL is the ratio between the square of the mean and the variance of a homogeneous region, where as Cx is given by the ratio between standard deviation and the mean intensity of a homogeneous region. The results in terms of ENL and Cx are given in Table 2. Our proposed method resulted in the highest ENL values in all 4 regions which signifies the best despeckling performance out of the considered approaches. Low Cx values indicate better preservation of texture and our proposed algorithm gave the lowest Cx values in all cases. From Fig. 2, we can observe that ENL and Cx results are consistent with the visual results.

## 4. CONCLUSION

We proposed a network architecture that encompasses a transformer-based encoder and a convolution-based decoder, for SAR image despeckling. When compared with several

Real SAR image      PPB      SAR-BM3D      SAR-CNN      ID-CNN      Proposed

**Fig. 2**. Results on real SAR images.

existing filter-based and CNN-based despeckling methods, the results on synthetic and real SAR images show promising quantitative and qualitative improvements. The proposed method proved to be effective in reducing speckle while preserving the texture and fine details in real SAR images.

## 5. REFERENCES

[1] Jong-Sen Lee, "Speckle analysis and smoothing of synthetic aperture radar images," *Computer Graphics and Image Processing*, vol. 17, no. 1, pp. 24–32, 1981.

[2] Darwin T. Kuan, Alexander A. Sawchuk, Timothy C. Strand, and Pierre Chavel, "Adaptive noise smoothing filter for images with signal-dependent noise," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PAMI-7, no. 2, pp. 165–177, 1985.

[3] Charles-Alban Deledalle, Loïc Denis, and Florence Tupin, "Iterative weighted maximum likelihood denoising with probabilistic patch-based weights," *IEEE Transactions on Image Processing*, vol. 18, no. 12, pp. 2661–2672, 2009.

[4] Sara Parrilli, Mariana Poderico, Cesario Vincenzo Angelino, and Luisa Verdoliva, "A nonlocal sar image denoising algorithm based on llmmse wavelet shrinkage," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 50, no. 2, pp. 606–616, 2012.

[5] Giulia Fracastoro, Enrico Magli, Giovanni Poggi, Giuseppe Scarpa, Diego Valsesia, and Luisa Verdoliva, "Deep learning methods for synthetic aperture radar image despeckling: An overview of trends and perspectives," *IEEE Geoscience and Remote Sensing Magazine*, vol. 9, no. 2, pp. 29–51, 2021.

[6] G. Chierchia, D. Cozzolino, G. Poggi, and L. Verdoliva, "Sar image despeckling through convolutional neural networks," in *2017 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, 2017, pp. 5438–5441.

[7] Puyang Wang, He Zhang, and Vishal M. Patel, "Sar image despeckling using a convolutional neural network," *IEEE Signal Processing Letters*, vol. 24, no. 12, pp. 1763–1767, 2017.

[8] Sergio Vitale, Giampaolo Ferraioli, and Vito Pascazio, "A new ratio image based cnn algorithm for sar despeckling," in *IGARSS 2019 - 2019 IEEE International Geoscience and Remote Sensing Symposium*, 2019, pp. 9494–9497.

[9] Shuaiqi Liu, Yu Lei, Luyao Zhang, Bing Li, Weiming Hu, and Yu-Dong Zhang, "Mrddanet: A multiscale residual dense dual attention network for sar image denoising," *IEEE Transactions on Geoscience and Remote Sensing*, pp. 1–13, 2021.

[10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," *ICLR*, 2021.

[11] Zhendong Wang, Xiaodong Cun, Jianmin Bao, and Jianzhuang Liu, "Uformer: A general u-shaped transformer for image restoration," *arXiv preprint 2106.03106*, 2021.

[12] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte, "Swinir: Image restoration using swin transformer," in *IEEE International Conference on Computer Vision Workshops*, 2021.

[13] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo, "Segformer: Simple and efficient design for semantic segmentation with transformers," *arXiv preprint arXiv:2105.15203*, 2021.

[14] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 568–578.

[15] D. Martin, C. Fowlkes, D. Tal, and J. Malik, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," in *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, 2001, vol. 2, pp. 416–423 vol.2.

[16] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, 2017.