

A CNN-Based Coherence-Driven Approach for InSAR Phase Unwrapping

Francescopaolo Sica[✉], Member, IEEE, Francesco Calvanese,
Giuseppe Scarpa[✉], Senior Member, IEEE, and Paola Rizzoli[✉]

Abstract—Phase unwrapping (PU) is among the most critical tasks in synthetic aperture radar (SAR) interferometry (InSAR). Due to the presence of noise, the interferogram usually presents phase inconsistencies, also called residues, which imply a nonunivocal solution. This work investigates the PU problem from a semantic segmentation perspective by exploiting convolutional neural network (CNN) models. In particular, by exploiting a popular deep-learning architecture, we introduce the interferometric coherence as an input feature and analyze the performance increase against classical methods. For the network training, we generate a variegated data set by introducing a controlled number of phase residues, and considering both synthetic and real InSAR data. Eventually, we compare the proposed method to state-of-the-art algorithms on synthetic and real InSAR data taken from the TanDEM-X mission, obtaining encouraging results.

Index Terms—Convolutional neural networks (CNNs), phase unwrapping (PU), SAR interferometry, semantic segmentation, synthetic aperture radar (SAR).

I. INTRODUCTION

SYNTHETIC aperture radar interferometry (InSAR) is one of the most applied techniques in remote sensing for the retrieval of ground topography and deformations. The InSAR technique exploits the phase difference between a pair of coregistered SAR images, namely the interferometric phase ψ . The bidimensional image of ψ is called interferogram and is typically measured modulo 2π . To reconstruct topographic heights or ground deformations from the interferogram, it is necessary to convert the wrapped interferometric phase into an absolute phase field ϕ , by adding the correct multiple of 2π to each fringe ($2\pi k$ with k the wrap count).

Such an ill-posed problem, known as phase unwrapping (PU), usually requires additional constraints to be properly addressed. The usual assumption is the so-called Itoh's condition [1], which restricts the phase difference between adjacent pixels to the range $[-\pi, \pi]$. Consequently, we can solve the PU problem by applying a univocal integration process along any spatial path. However, the interferogram may present

Manuscript received July 10, 2020; revised September 3, 2020 and October 1, 2020; accepted October 3, 2020. Date of publication October 21, 2020; date of current version December 21, 2021. (Corresponding author: Francescopaolo Sica.)

Francescopaolo Sica and Paola Rizzoli are with the German Aerospace Center, Microwaves and Radar Institute, 82234 Wessling, Germany (e-mail: francescopaolo.sica@dlr.de; paola.rizzoli@dlr.de).

Francesco Calvanese and Giuseppe Scarpa are with the Department of Electrical Engineering and Information Technology, University of Napoli Federico II, 80138 Napoli, Italy (e-mail: fran.calvanese@studenti.unina.it; giuseppe.scarpa@unina.it).

Digital Object Identifier 10.1109/LGRS.2020.3029565

phase jumps greater than π , which may occur because of the true signal structure or, more commonly, in the presence of noise. The latter are phase inconsistencies, called residues, which impair the integration procedure.

Several approaches have been proposed in the literature to solve the PU problem. Path following algorithms, such as the Branch-Cut (BC) method [2], exploit the knowledge on residues' position to select the best integration path to unwrap the phase. This kind of approach is considered reliable in the presence of low residues' density. Differently, other methods minimize the L^1 and L^2 norms of the difference between unwrapped and wrapped phase gradients. On the one hand, the L^1 norm is used in network programming solutions for loss minimization as in [3] and [4]. Approaches based on these methods are very efficient and are extensively applied still today. On the other hand, the L^2 norm is used in the least-squares (LS) methods, [5] and [6], which are based on the assumption of a smooth phase field. Besides, the PU problem can be addressed in a Bayesian framework, combining *a priori* model on ϕ [7]. Representative solutions of this last category are Statistical-Cost, Network-Flow Algorithm for Phase Unwrapping (SNAPHU) [8] and PU via MAX flows (PUMA) [9]. Recently, deep learning methods for PU have appeared in the literature as well. To the best of our knowledge the first attempt has been proposed in [10], where the PU problem was converted into a segmentation task. In Yan *et al.* [11], the authors embedded phase denoising and wrap count reconstruction in a single framework. A similar approach for InSAR unwrapping has been proposed in [12] with the difference that instead of absolute wrap count values, its gradient is reconstructed. This approach improves the algorithm's generality, which would otherwise be impaired depending on the wrap count dynamic. However, none of the mentioned works consider the interferometric coherence as an additional source of information to support the PU task.

In this work, we propose a new PU method based on convolutional neural networks (CNNs). The objective is to exploit the intrinsic capability of CNNs to solve spatial-dependent problems, to better handle phase inconsistencies that tend to remain unsolved with traditional algorithms. We approach the problem as a semantic segmentation task, including the coherence as an input feature to the network, to enforce phase consistency and achieve improved local accuracy.

The rest of the letter is organized as follows. Section II presents the proposed method. Section III describes the experimental results on synthetic and real data, while conclusions and outlook are finally drawn in Section IV.

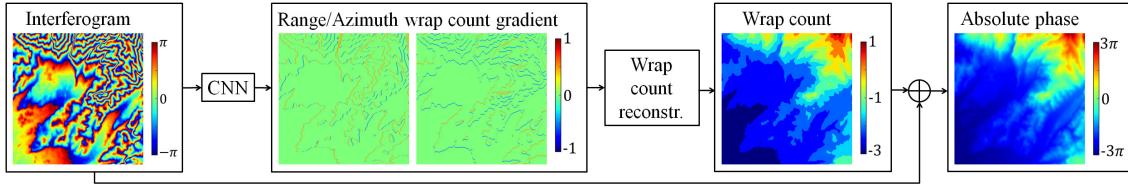


Fig. 1. Phase field. High-level flowchart of the proposed concept.

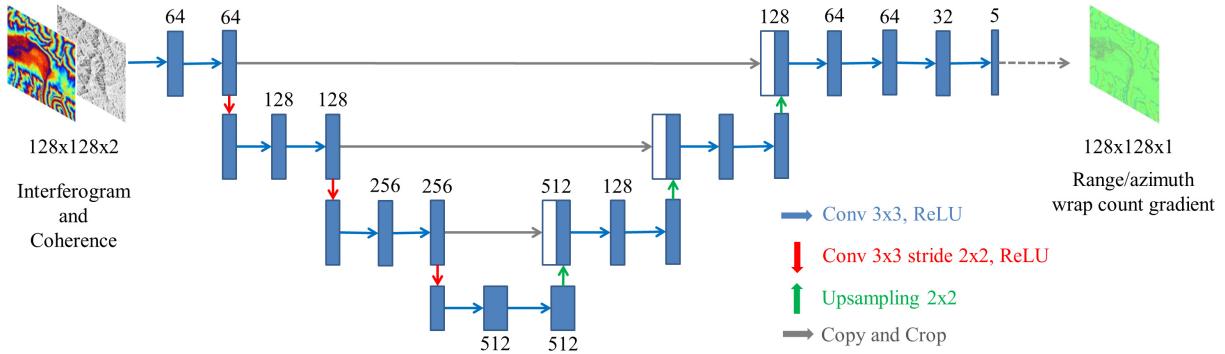


Fig. 2. Proposed CNN model, based on the original U-Net architecture.

II. PROPOSED METHOD

Let us consider the interferometric signal model used in [13]. The observed complex interferogram γ can be written as

$$\gamma = A\rho e^{j\psi} + w \quad (1)$$

where w is a zero-mean noise contribution that depends on the following three parameters: the SAR amplitude A , the interferometric phase ψ , and the coherence ρ . The latter represents the correlation coefficient between the interferometric pair and is the key quantity for assessing the quality of an interferogram. For each image location s , the relationship between wrapped (ψ) and unwrapped (ϕ) phase is simply given by

$$\phi(s) = \psi(s) + 2\pi k(s) \quad (2)$$

being $k(s)$ the wrap count at point s to be estimated.

To secure a robust behavior of the proposed model, we prefer to formulate the wrap count estimation in differential terms by targeting the variations of $k(s)$, rather than its absolute values, and then properly integrating them. This choice helps the network during the training phase by making it insensitive to the absolute phase value and eventually allowing for a better generalization capability. Specifically, the wrap count derivatives are independently predicted for range and azimuth directions, leaving to the subsequent integration process the aim of preserving the consistency of the provided 2-D maps. Therefore, we refer for simplicity to either the range or the azimuth case, since we can straightforwardly extend the same model to the other case. We consider the following approximations of the wrap count derivative $\Delta_k(s)$ along the

linear coordinate s :

$$\Delta_k(s) \triangleq \begin{cases} k(s+1) - k(s), & s \text{ leftmost} \\ k(s) - k(s-1), & s \text{ rightmost} \\ (k(s+1) - k(s-1))/2, & \text{otherwise} \end{cases}. \quad (3)$$

Moreover, by assuming that no jumps greater than 2π occur between adjacent pixels ($k(s) - k(s-1) \in \{0, \pm 1\}$), it follows that $\Delta_k(s)$ can assume only five values:

$$\Delta_k(s) \in \{0, \pm 1/2, \pm 1\}.$$

Under these assumptions, the prediction of $\Delta_k(s)$ can be regarded as a semantic segmentation problem with five possible classes that we predict with two CNNs (detailed below), for the range and azimuth directions, $\Delta_k^r(s)$ and $\Delta_k^a(s)$, respectively. Finally, the wrap count derivatives are eventually integrated to reconstruct the 2-D wrap count map, which has to be added to the original interferogram to retrieve the absolute phase field. Fig. 1 summarizes the main developed concept.

A. CNN Architecture Design

We base our network architecture on the U-Net [14]. This choice stems from the observation that this network, designed initially for semantic segmentation of medical images, has proven its effectiveness in many different applications. Here, we modify the layers parameters, the network depth, and the final layers after the decoder stage, to specifically solve the PU problem. Our network architecture is depicted in Fig. 2. It consists of an encoder section on the left side and a decoder one on the right side, resulting in a U-like structure. The encoder path follows the typical architecture of a CNN for classification. Indeed several stages associated with different scales can be identified. Each stage consists of a repeated application of two 3×3 convolutions, interleaved by rectified linear unit (ReLU)

activations, with a terminal 3×3 convolutional layer with a 2×2 stride for the downsampling, which doubles the number of features. By doing so, the next stage works on a deeper but reduced-scale feature volume. Such a depth is kept constant along the stage till the last (strided) convolution that is responsible for moving to the next stage (see Fig. 2 to read the number of features flowing in each stage). Moreover, at the beginning of each stage, a batch normalization (BN) layer is applied to speed up the network training and mitigate its dependence on the initialization [15]. The decoder section comprises dual stages that allow for the progressive restoration of the retrieved feature's spatial resolution, simultaneously reducing their number, by halving it at each 2×2 upscaling (up-convolution). Therefore, the upsampled features are concatenated with the corresponding (same scale) features coming from the encoding path, thanks to skip connections. Concatenated feature blocks are then convolved twice using ReLU activations, similarly as in the encoder. The exit decoding stage comprises two additional convolutional layers which compact the output in five features provided with a five-way softmax activation, finally representing the desired pixel-wise class membership probabilities. The network is fed with both the interferogram and the coherence, while the output is the range or the azimuth component of the wrap count gradient. The use of the coherence as additional input provides the network with a local indication about the presence of residues, therefore, supporting the correct processing of critical noisy regions of the interferogram. Finally, it is worth remarking that, although the overall architecture is identical for the range and the azimuth cases, disjoint training procedures were carried out for each of them, eventually leading to different network parameters. The choice of operating the CNN model separately to the range and azimuth directions resides in the fact that SAR geometry is highly nonisotropic due to the side-looking nature of SAR. Experimental results have confirmed this peculiarity through performance improvement. We don't show these experiments for the sake of brevity.

B. Generation of the Training Data Set

To properly train our network, ensuring a good generalization capability, we generate a hybrid input data set that relies on synthetic and real InSAR data. We first simulate interferometric phase patterns by back-geocoding the edited Shuttle Radar Topography Mission (SRTM) digital elevation model (DEM) [16]. As reference SAR geometries, we utilize the acquisition parameters of ten TanDEM-X [17] Stripmap sample acquisitions over the Austrian Alps, with orthogonal baselines between 50 and 140 m. In this way, we can assure the illumination of a variegated topography (from flat to high-relief terrain) and different land cover classes, such as vegetation, bare soil, and agricultural areas. Referring now to the signal model of (1), we generate three different training data sets by utilizing the synthetic phase ψ from SRTM and varying the amplitude and coherence as follows:

Case 1: noise-free model, where $A = 80$ and $\rho = 1$ (considering an overall dynamic range from 0 up to 255 for A and from 0 and 1 for ρ). We consider this case to support

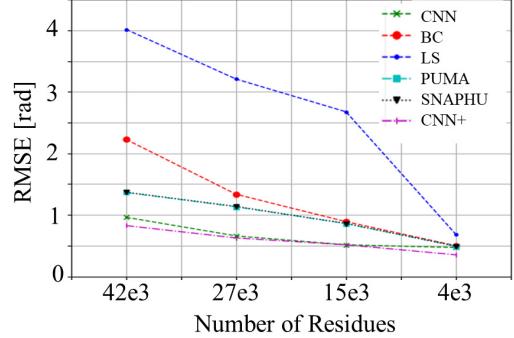


Fig. 3. Mean RMSE between the predicted unwrapped phase and the reference noise-free one for our proposed CNN (CNN+), CNN without coherence (CNN), BC, LS, PUMA, and SNAPHU, with respect to the total number of residues.

the network in understanding the general PU problem, without introducing additional signal corruption sources.

Case 2: $A = 80$, while the coherence ρ is constant all over the patch and can assume values between 0.25 and 0.95. Both interferometric phase and coherence are then filtered using a boxcar window of size 3×3 pixels. In this way, we simulate the presence of residual noise after phase estimation.

Case 3: the amplitude A and the coherence ρ are derived from TanDEM-X real InSAR data, acquired over the considered region and with the same acquisition geometry used for back-geocoding the SRTM DEM, using a standard 3×3 averaging window.

C. Training Process

We trained the network by employing the *Adam* algorithm (initial learning rate of 10^{-4}). Each minibatch for training comprises 32 input–output examples of size 128×128 , randomly sampled from an overall training-validation set comprising 91 607 samples (80% reserved for training). The training comprises two consecutive steps. First the network is trained using the simulated wrapped phase patches described at *Case 1* and *Case 2* of Section II-B and using the corresponding noise-free horizontal (range) and vertical (azimuth) components of the wrap count derivative as output reference. Afterward, a second training is performed for fine-tuning by exploiting the data of *Case 3* and using the same optimizer and mini-batch size but with an initial learning rate of 10^{-6} . Overall, we allocate 10% of the entire training data set for *Case 1*, 50% for *Case 2*, and 40% for *Case 3*. The small percentage of *Case 1* is because (being noise-free) is easy to learn by the network. Therefore, a lower number of examples are sufficient. Besides, the percentage of *Case 2* and *Case 3* was experimentally chosen, allowing for a good balance between synthetic and real training data. We trained on a NVIDIA Titan X with 12 GB of GPU memory for 150 epochs for the first and 50 for the second steps.

Moreover, as loss function \mathcal{L}_{Tot} , we use a combination of three different components

$$\mathcal{L}_{\text{Tot}} = \mathcal{L}_{\text{cross}} + \mathcal{L}_{L_1} + \mathcal{L}_{\text{jacc}} \quad (4)$$

where $\mathcal{L}_{\text{cross}}$ is the cross-entropy loss, applied pixel-wise for each predicted class label at the output of the classifier.

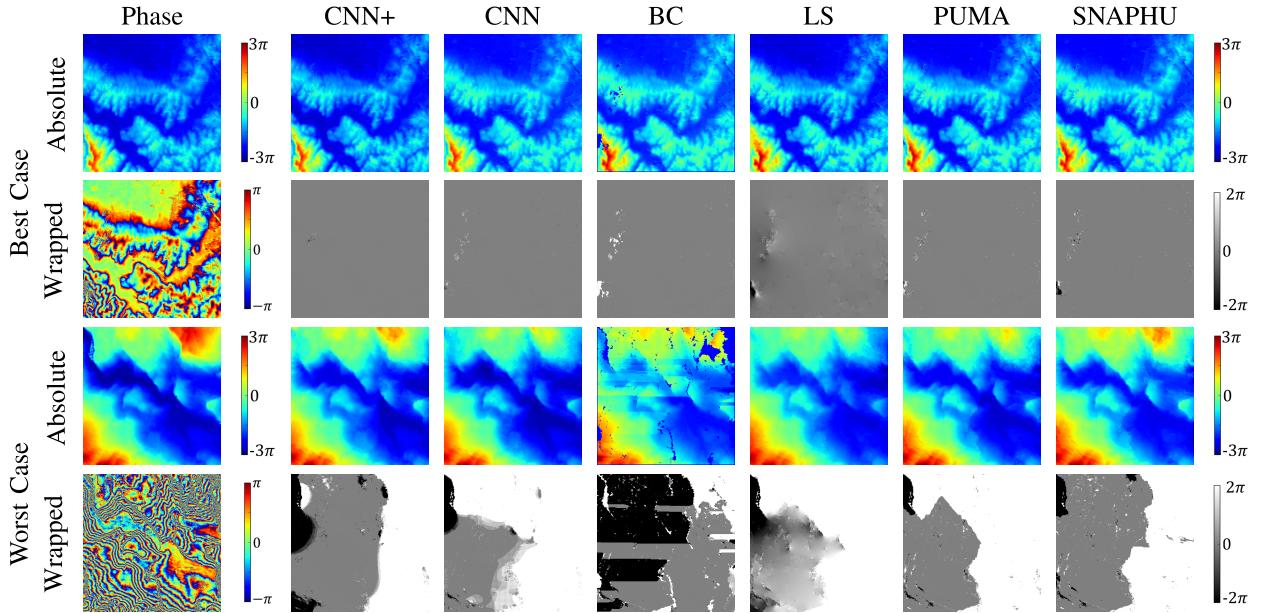


Fig. 4. Absolute phase from SRTM and TanDEM-X wrapped interferogram (first column), unwrapped phases (first and third rows), and error images (second and fourth rows) for the best and worst case scenarios, respectively.

\mathcal{L}_{L_1} is the L^1 norm between the obtained prediction and the ground-truth, which we have experimentally found to speed-up the training process [18]. Eventually, \mathcal{L}_{jacc} is the Jaccard distance loss, which provides a gradient-consistent and accurate estimation of the wrap count derivatives and, therefore, a precise prediction of the absolute unwrapped phase. Finally, as already mentioned in Section II-A, we separately train the described network architecture to estimate the azimuth and range derivatives of the wrap count.

D. Wrap Count Reconstruction

Once the $n \times n$ range and azimuth components of the wrap count gradient, Δ_k^r and Δ_k^a , respectively, are estimated, we must reconstruct the complete 2-D image of the wrap count k [19]. To do so, we first vectorize each gradient matrix Δ_k^r and Δ_k^a into $n^2 \times 1$ column vectors, f_r and f_a , respectively, that are concatenated in a single $2n^2 \times 1$ vector $f \triangleq [f_r^T, f_a^T]^T$. In this way, the overall differential operation for a $n \times n$ sample reduces to a simple matrix multiplication between the vectorized wrap count, say k_v , and the derivative matrix Γ of dimension $2n^2 \times n^2$ that contains the derivation coefficients: $f = \Gamma k_v$. This overcomplete equation system can be easily solved with respect to k_v using the pseudo-inverse matrix of Γ

$$k_v = (\Gamma^T \Gamma)^{-1} \Gamma^T f. \quad (5)$$

Finally, the predicted wrap count map k is given by the integer approximation of the $n \times n$ -reshaped version of k_v .

III. EXPERIMENTAL RESULTS

In this section, we present two experiments for evaluating the performance of the proposed methodology. On the one hand, we demonstrate the robustness of the implemented algorithm with respect to different noise levels (and therefore, to residues density) in a controlled environment by exploiting

synthetic data. On the other hand, we evaluate the performance on real InSAR data using TanDEM-X acquisitions. In both cases, we compare the performance with the following state-of-the-art PU methods: the BC [2], LS [5], SNAPHU [8], and PUMA [9]. Additionally, to further quantify the impact of using the coherence as an input feature to the network, we also consider a baseline case in which the noisy interferogram is the only input feature to the proposed CNN architecture. For the first experiment, we consider the synthetic test case presented in Fig. 1, which was not previously used during the training phase. Starting from the noise-free interferogram, we generate a series of noisy interferograms by introducing a noise component in the signal model, as presented in Section II. We utilize patches of 128×128 pixels and different constant images of coherence, with a variation comprised between 0.65 and 0.95 and a regular increment of 0.1. It is noted that no denoising procedure is applied in this case. This choice relies on the fact that we aim to analyze the proposed methods in the presence of an increasing number of residues in the data, without adding further error sources. As a quality measure, we then compute the root mean square error (RMSE) between the predicted unwrapped phase and the reference noise-free one from Fig. 1. The results are summarized in Fig. 3. The proposed CNN-based method (CNN+) has the lowest RMSE for each considered case. As expected, all algorithms show similar performance in the presence of a low density of residues. On the contrary, it is worth noting how the CNN performs better at a high density of residues thanks to its capability to solve spatial-dependent problems. For the second experiment, we now consider a set of ten different patches of 512×512 pixels from a real TanDEM-X single-pass interferogram. For each patch and each considered algorithm, we compute the error image between the estimated unwrapped phase field and the reference absolute phase, obtained by back-geocoding the SRTM DEM. The results are summarized in Table I, which displays the

TABLE I

RMSE [RADIAN] FOR THE CONSIDERED METHODS. **OVERALL**: AVERAGE AND STANDARD DEVIATION (IN BRACKETS) OF THE RMSE OVER THE TEN CONSIDERED TEST PATCHES. **BEST** AND **WORST**: RMSE FOR THE BEST AND WORST-CASE TEST PATCHES, RESPECTIVELY

PU Method	RMSE					
	CNN+	CNN	BC	LS	PUMA	SNAPHU
Overall	3.90 (2.72)	4.52 (3.5)	10.05 (5.75)	8.77 (5.83)	4.62 (3.68)	4.08 (3.02)
Best	0.11	0.26	1.06	0.62	0.34	0.41
Worst	8.64	12.65	18.20	17.89	11.98	9.07

average RMSE and standard deviation of RMSE values over the ten considered test patches, together with the RMSE for the best and worst-case test patches. For visual inspection, the corresponding unwrapped phase and error images for the best and worst cases are depicted in Fig. 4, together with the images of the reference absolute phase derived from SRTM and of the input interferogram. We observe that the proposed method (CNN+) shows the best performance with respect to all state-of-the-art algorithms, resulting in an overall lower mean RMSE. Specifically, if we consider the best case scenario, CNN+ shows a further improvement with respect to the baseline CNN case. This behavior is probably related to the spatial-dependent nature of the PU problem. Indeed, CNN+ is able to learn additional information about the local noise distribution through the utilization of the coherence map and, therefore, to better handle the presence of high-density residues. Consistently, the RMSE considerably increases in the worst-case scenario, where a low average coherence makes PU fail over a considerable portion of the patch. In particular, so-called 2π -ambiguities arise in correspondence of adjacent regions separated by very low coherence areas, and sudden phase jumps of 2π appear. Even though none of the considered methods can fully avoid this phenomenon, CNN+ shows the best performance. Together with SNAPHU, they are the two algorithms that better cope with the problem, reducing the extension of the affected areas.

IV. CONCLUSION AND OUTLOOK

In this letter, we have introduced a novel method based on a CNN model that estimates the wrap count gradient to derive the unwrapped phase field. The interferometric coherence plays a twofold key role in the whole process. On the one hand, it drives the characterization of noise during the definition of the training data set. On the other hand, when used as an additional input feature, it helps the network identify and manage critical noisy regions. From the presented validation, we can conclude that the proposed method is robust and very promising also in the presence of high-density residues, showing state-of-the-art performance. The use of the coherence feature results in overall performance improvement on both synthetic and real data. As future developments, we plan to refine the training data set and fine-tune the CNN parameters for other sensors and applications. Moreover, it will be our priority to compare the performance to other state-of-the-art deep learning-based methods as well, such as

the one presented in [12], when the trained model will be released. It could also be of interest to explore different approaches that abstract from the use of input engineered features, for example, by letting the coherence be directly estimated from the CNN model. This approach has already been investigated in the literature, showing a high potential when applied in a supervised manner [20]. Moreover, as shown in [21], an unsupervised learning approach could also help to overcome training problems caused by the scarcity of data, which is a common issue in remote-sensing applications.

REFERENCES

- [1] K. Itoh, "Analysis of the phase unwrapping algorithm," *Appl. Opt.*, vol. 21, no. 14, p. 2470, 1982.
- [2] R. M. Goldstein, H. A. Zebker, and C. L. Werner, "Satellite radar interferometry: Two-dimensional phase unwrapping," *Radio Sci.*, vol. 23, no. 4, pp. 713–720, Jul. 1988.
- [3] T. J. Flynn, "Two-dimensional phase unwrapping with minimum weighted discontinuity," *JOSA A*, vol. 14, no. 10, pp. 2692–2701, 1997.
- [4] M. Costantini, "A novel phase unwrapping method based on network programming," *IEEE Trans. Geosci. Remote Sens.*, vol. 36, no. 3, pp. 813–821, May 1998.
- [5] D. L. Fried, "Least-square fitting a wave-front distortion estimate to an array of phase-difference measurements," *JOSA*, vol. 67, no. 3, pp. 370–375, 1977.
- [6] H. A. Zebker and Y. Lu, "Phase unwrapping algorithms for radar interferometry: Residue-cut, least-squares, and synthesis algorithms," *JOSA*, vol. 15, no. 3, pp. 586–598, 1998.
- [7] A. Budillon, G. Ferraiuolo, V. Pascazio, and G. Schirinzi, "Multichannel SAR interferometry via classical and Bayesian estimation techniques," *EURASIP J. Adv. Signal Process.*, vol. 2005, no. 20, p. 3180, Dec. 2005.
- [8] C. W. Chen and H. A. Zebker, "Phase unwrapping for large SAR interferograms: Statistical segmentation and generalized network models," *IEEE Trans. Geosci. Remote Sens.*, vol. 40, no. 8, pp. 1709–1719, Aug. 2002.
- [9] J. M. Bioucas-Dias and G. Valadão, "Phase unwrapping via graph cuts," *IEEE Trans. Image Process.*, vol. 16, no. 3, pp. 698–709, Mar. 2007.
- [10] T. Zhang *et al.*, "Rapid and robust two-dimensional phase unwrapping via deep learning," *Opt. Express*, vol. 27, no. 16, p. 23173, Aug. 2019.
- [11] K. Yan, Y. Yu, T. Sun, A. Asundi, and Q. Kemao, "Wrapped phase denoising using convolutional neural networks," *Opt. Lasers Eng.*, vol. 128, pp. 1–5, May 2020.
- [12] L. Zhou, H. Yu, and Y. Lan, "Deep convolutional neural network-based robust phase gradient estimation for two-dimensional phase unwrapping using SAR interferograms," *IEEE TGRS*, vol. 58, no. 7, pp. 4653–4665, Jul. 2020.
- [13] F. Sica, D. Cozzolino, L. Verdoliva, and G. Poggi, "The offset-compensated nonlocal filtering of interferometric phase," *Remote Sens.*, vol. 10, no. 9, p. 1359, Aug. 2018.
- [14] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Proc. Int. Conf. Med. Image Comput. Comput.-Assist. Intervent.* Cham, Switzerland: Springer, 2015.
- [15] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proc. Int. Conf. Mach. Learn.*, vol. 37, 2015, pp. 1–4.
- [16] T. G. Farr *et al.*, "The shuttle radar topography mission," *Rev. Geophys.*, vol. 45, pp. 1–33, May 2007.
- [17] G. Krieger *et al.*, "TanDEM-X: A satellite formation for high-resolution SAR interferometry," *IEEE Trans. Geosci. Remote Sens.*, vol. 45, no. 11, pp. 3317–3341, Nov. 2007.
- [18] G. Scarpa, S. Vitale, and D. Cozzolino, "Target-adaptive CNN-based pansharpening," *IEEE Trans. Geosci. Remote Sens.*, vol. 56, no. 9, pp. 5443–5457, Sep. 2018.
- [19] G. Farneback, J. Rydell, T. Ebbers, M. Andersson, and H. Knutsson, "Efficient computation of the inverse gradient on irregular domains," in *Proc. IEEE 11th Int. Conf. Comput. Vis.*, 2007, pp. 1–8.
- [20] F. Sica, G. Gobbi, P. Rizzoli, and L. Bruzzone, " Φ -net: Deep residual learning for InSAR parameters estimation," *IEEE Trans. Geosci. Remote Sens.*, pp. 1–25, 2020.
- [21] S. Mukherjee, A. Zimmer, X. Sun, P. Ghuman, and I. Cheng, "An unsupervised generative neural network approach for InSAR phase filtering and coherence estimation," *IEEE Geosci. Remote Sens. Lett.*, early access, Jul. 30, 2020, doi: [10.1109/LGRS.2020.3010504](https://doi.org/10.1109/LGRS.2020.3010504).