

# SEA ICE SEGMENTATION FROM SAR DATA BY CONVOLUTIONAL TRANSFORMER NETWORKS

Nicolae-Cătălin Ristea<sup>1</sup>, Andrei Anghel<sup>1</sup>, Mihai Datcu<sup>1,2</sup>

CEOSpaceTech, University Politehnica of Bucharest, Romania<sup>1</sup>  
Remote Sensing Technology Institute, German Aerospace Center (DLR), Germany<sup>2</sup>

## ABSTRACT

Sea ice is a crucial component of the Earth’s climate system and is highly sensitive to changes in temperature and atmospheric conditions. Accurate and timely measurement of sea ice parameters is important for understanding and predicting the impacts of climate change. Nevertheless, the amount of satellite data acquired over ice areas is huge, making the subjective measurements ineffective. Therefore, automated algorithms must be used in order to fully exploit the continuous data feeds coming from satellites. In this paper, we present a novel approach for sea ice segmentation based on SAR satellite imagery using hybrid convolutional transformer (ConvTr) networks. We show that our approach outperforms classical convolutional networks, while being considerably more efficient than pure transformer models. ConvTr obtained a mean intersection over union (mIoU) of 63.68% on the AI4Arctic data set, assuming an inference time of 120ms for a 400 × 400 km<sup>2</sup> product.

**Index Terms**— transformers, remote sensing, SAR, deep learning, semantic segmentation.

## 1. INTRODUCTION

Sea ice retreat, particularly in the Arctic, has been one of the most significant responses to global climate change. Therefore, sea ice cover and sea ice concentration are vital parameters for conducting climate change research and navigation in polar regions. To support the logistics for the transport industry, there is a high demand for local-scale high-resolution information on Arctic marine conditions. Such information is critical for operations planning, shipping routes and sustainable development of the North [1]. To this end, the research infrastructure in the ice covered areas have grown significantly in the last decades. Satellite based synthetic aperture radar (SAR) systems have been employed to monitor the vast regions of the Arctic (e.g., RADARSAT-2, RADARSAT Constellation Mission, Sentinel-1A and -1B). These systems have a high spatial resolution and have regional coverage (e.g., up

to 500 km by 500 km), making them ideal for monitoring large regions. Nevertheless, even if the satellite infrastructure assures a high amount of data, this needs to be processed and interpreted to extract key information. Considering the amount of data, manual processing is barely possible, therefore automated algorithms are needed.

Recently, deep learning models have been widely adopted for the geoscience field [2–7]. Those models have the potential to greatly improve the efficiency and accuracy of sea ice analysis, as well as to enable the analysis of large amounts of data that would be impractical to process manually. In [4] the authors propose a U-Net architecture for sea ice segmentation on data acquired from Sentinel-1 and manually labelled. The authors observed that even if the network was trained with a limited amount of products, the network is still able to learn the patterns and obtain high segmentation scores. Gao et al. [5] employed a dense neural architecture to detect changes over time in the sea ice areas. They adopted a transfer learning strategy to increase the network performance. Closer to our work, in [6] the authors proposed a semi-supervised algorithm based on graph convolutional networks for sea ice segmentation. They obtained superior results compared with a ResNet based architecture in the limited data scenario.

Lately, inspired by the success of self-attention layers and transformer architectures in the computer vision field [8, 9], there have been employed transformer architectures in the geoscience field [10–14] with remarkable results. For example, in [10] the authors use a SwinTransformer architecture, originally employed for computer vision tasks, to detect the melt ponds on Arctic sea ice. They developed a cross-channel attention into the decoder block, which boosts the model’s performance.

Distinct from all mentioned methods, we propose a hybrid convolutional transformer (ConvTr) model, which combines the benefits of convolutional networks (e.g., efficiency) and transformer blocks (e.g., global attention). The architecture has a transformer core, designed to compute attention based features in a smaller latent space. Moreover, our approach uses a large scale SAR data set, AI4Arctic [15], enabling the network to learn general patterns from diverse scenes, instead of overfitting the training on a limited number of products.

In summary, our contribution is twofold:

This work was supported by the grant of the Romanian Ministry of Education and Research, CNCS - UEFISCDI, project number PN-III-P4-ID-PCE-2020-2120, within PNCDI III.

- We propose a hybrid convolutional transformer architecture for sea ice segmentation, obtaining the best time-accuracy trade-off.
- We train our model on a large scale data set, showing the generalisation capacity of our network.

## 2. METHOD

### 2.1. Data preprocessing

Considering that the products from the AI4Arctic [15] data set have a high dimensional size, being impractical to train neural networks on the entire dimension, we crop fix length windows of size  $P \times P$ , where  $P \in \mathbb{N}$  from the product. The cropping is performed such that, in each window there is more than a single label from all considered (*sea, ice, land*). A training sample of size  $2 \times P \times P$  contains both HH and HV polarizations. Finally, each sample is normalised before being ready to be fed into the network.

### 2.2. Convolutional Transformer architecture

In our work, we employed the ConvTr architecture composed of a convolutional downsampling block, a convolutional transformer block, and a deconvolutional upsampling block, as illustrated in Fig. 1. We highlight that, without the convolutional downsampling block and the replacement of dense layers with convolutional layers inside the transformer block [16], the transformer would not be able to learn to segment images larger than  $128 \times 128$  pixels, due to memory overflow (measured on a Nvidia GeForce RTX 3090 GPU with 24GB of VRAM).

**Downsampling block.** The downsampling block starts with a convolutional layer formed of 32 filters with a spatial support of  $7 \times 7$ , which are applied using a padding of 3 pixels to preserve the spatial dimension, while enriching the number of feature maps to 32. Next, we apply three convolutional layers composed of 32, 64 and 128 filters, respectively. All convolutional filters have a spatial support of  $3 \times 3$  and are applied at a stride of 2, using a padding of 1. Each layer is followed by batch-norm [17] and Rectified Linear Units (ReLU) [18].

**Transformer block.** The downsampling block is followed by the convolutional transformer block, which preserves the spatial dimension between the input and output tensors. The convolutional transformer block is inspired by the block proposed in [2]. More precisely, the input tensor is interpreted as a set of overlapping tokens (patches from the input tensor). The sequence of tokens is projected onto a set of weight matrices implemented as depth-wise separable convolution operations. The convolutional projection is formed of three nearly identical projection blocks, with separate parameters. The output query, keys and values ( $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ ) are passed to a multi-head attention layer, with the goal of capturing the interaction among all tokens by encoding each entity in terms

of the global contextual information. Next, the output passes through a batch-norm and a pointwise convolution, with the corresponding residual connections. The process is repeated  $L$  times, which denotes the depth of the transformer block. The block is visually described in Fig. 1.

**Upsampling block.** The last block of our ConvTr applies upsampling operations, being designed to revert the transformation of the downsampling block. The upsampling block is formed of three transposed convolutional layers comprising 128, 64 and 32 filters, respectively. All kernels have a spatial support of  $3 \times 3$ , being applied at a stride of 2, using a padding of 1. Similar to the downsampling block, we apply batch-norm and ReLU activations after each transposed convolutional layer. Finally, we employ a convolutional layer with  $C \in \mathbb{N}$  filters, each filter having a kernel dimension of  $7 \times 7$  and a padding of 3, to reduce the number of channels from 32 to  $C$ . In this manner, we obtained the same output dimension as the input image, where each of the  $C$  channels represents the probability that a certain pixel is part of the corresponding class.

### 2.3. Loss function

Considering the imbalanced classes from the AI4Arctic [15] data set, we optimised the model in accordance with the focal loss function [19]. In this manner, the network converged faster and was more robust at testing time to the minority class. Formally, the loss is defined below.

$$\mathcal{L} = \alpha(1 - p_t)^\gamma \mathcal{L}_{CE}, \quad (1)$$

where  $\alpha$  control the class weights,  $p_t$  is the probability of predicting the ground truth class,  $\gamma$  controls the degree of down-weighting of easy-to-classify pixels and  $\mathcal{L}_{CE}$  is the cross entropy loss function.

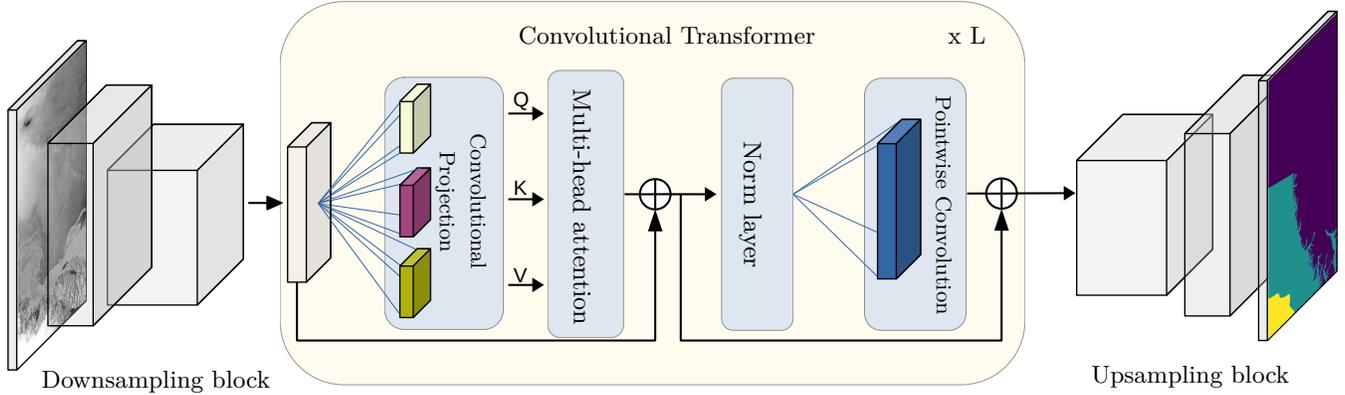
## 3. EXPERIMENTS

### 3.1. Data set

The AI4Arctic [15] Sea Ice data set are produced for the AI4EO sea ice competition initiated by the European Space Agency. The data set contains Sentinel-1 active microwave SAR data and corresponding passive MicroWave Radiometer data from the AMSR2 satellite sensor. Each product has associated ice charts that have been produced by the Greenland ice service at the Danish Meteorological Institute and the Canadian Ice Service for the safety of navigation. The scenes are from the time period from January 8 2018 to December 21 2021. The extra wide swath GRDM products cover a region of  $400 \times 400 \text{ km}^2$ , have a resolution of 90 meters and a pixel spacing of 40 meters. The entire data set contains 513 annotated products. We split the data into a training set (400 products) and a test set (113 products).

### 3.2. Hyper-parameters tuning

ConvTr is optimised with Adam using the focal loss function [19]. We start with an initial learning rate of  $10^{-4}$  and



**Fig. 1.** ConvTr segmentation architecture. The model is composed by a downsampling block comprising convolutional layers, a convolution transformer block comprising a multi-head self-attention mechanism, and an upsampling block comprising transposed convolutions.

**Table 1.** Segmentation and inference time results on the AI4Arctic [15] test set. ConvTr is compared against two baseline methods (ResNet AE, UNet [20]). We included for ablation the ConvTr only with convolutional blocks (AE) and only with transformer block (Transformer).

Method	mIoU (%)	Inference time (ms)
ResNet AE	53.04	87
UNet [4]	56.43	92
AutoEncoder	49.75	65
Transformer	63.81	473
ConvTr (ours)	63.68	120

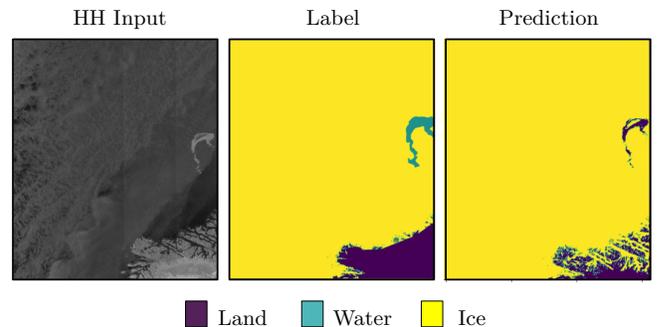
use a decay factor of 0.5 after every 10 epochs. We train each model for 50 epochs on mini-batches of 16 samples. We set the number of blocks to  $L = 5$  and each block has 5 attention heads. Regarding the training patch size, we found the optimal value to be  $P = 512$ .

### 3.3. Evaluation metrics

Since we perform semantic segmentation between three classes, we found the most insightful metric to be mean intersection over union (mIoU). The metric captures the overall performance of the models, regardless of the unbalanced distribution between classes. In addition, we reported the inference time for a full resolution scene, which has a spatial dimension about  $1100 \times 1100$ .

### 3.4. Results

In Table 1 we report the results for ConvTr against two baseline methods, ResNet based auto-encoder and UNet [20], used in [4]. We observe that ConvTr surpasses with more than 7% both baseline methods, while marginally raising the inference time. Regarding the importance of different parts of the network, when we only use the downsampling and



**Fig. 2.** Results obtained with ConvTr model for product 20180607T184326. Along with the prediction, we also included the HH input and the label.

upsampling blocks, we have the best inference speed, but the performance is drastically affected. If we employ only the transformer block, we note that the speed is highly impacted, while the accuracy is with only 0.13% higher. Therefore, combining both architectures, we exploit the benefits of those, attaining the best performance-speed trade-off.

In addition to the objective metrics, we included in Fig. 2 the result of our best ConvTr on 20180607T184326 product from the test set. We observe that the ice and land classes are well segmented, while the water class is miss classified. A potential reason could be the imbalanced training set.

## 4. CONCLUSION

In this paper, we propose a hybrid convolutional transformer architecture for sea ice segmentation, based on SAR data. We trained our model on a large scale data set, testing the generalisation capacity on over 100 products. Moreover, we showed that our hybrid architecture attains the best performance-speed trade-off, being feasible to be deployed for automated segmentation.

## 5. REFERENCES

- [1] Natalie Ann Carter, Jackie Dawson, Jenna Joyce, Annika Ogilvie, and Melissa Weber, “Arctic Corridors and Northern Voices: Governing Marine Transportation in the Canadian Arctic (Pond Inlet, Nunavut Community Report),” 2018.
- [2] Nicolae-Cătălin Ristea, Andrei Anghel, Mihai Datcu, and Bertrand Chapron, “Guided Unsupervised Learning by Subaperture Decomposition for Ocean SAR Image Retrieval,” *arXiv preprint arXiv:2209.15034*, 2022.
- [3] Nicolae-Cătălin Ristea, Andrei Anghel, Mihai Datcu, and Bertrand Chapron, “Guided Deep Learning by Subaperture Decomposition: Ocean Patterns from SAR Imagery,” in *Proceedings of IGARSS*. IEEE, 2022, pp. 6825–6828.
- [4] Yibin Ren, Huan Xu, Bin Liu, and Xiaofeng Li, “Sea Ice and Open Water Classification of SAR Images Using a Deep Learning Model,” in *Proceedings of IGARSS*. IEEE, 2020, pp. 3051–3054.
- [5] Yunhao Gao, Feng Gao, Junyu Dong, and Shengke Wang, “Transferred Deep Learning for Sea Ice Change Detection From Synthetic-Aperture Radar Images,” *IEEE Geoscience and Remote Sensing Letters*, vol. 16, no. 10, pp. 1655–1659, 2019.
- [6] Mingzhe Jiang, Xinwei Chen, Linlin Xu, and David A Clausi, “Semi-Supervised Sea Ice Classification of SAR Imagery Based on Graph Convolutional Network,” in *Proceedings of IGARSS*. IEEE, 2022, pp. 1031–1034.
- [7] Francesco Lattari, Borja Gonzalez Leon, Francesco Asaro, Alessio Rucci, Claudio Prati, and Matteo Matteucci, “Deep Learning for SAR Image Despeckling,” *Remote Sensing*, vol. 11, no. 13, pp. 1532, 2019.
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al., “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale,” in *Proceedings of ICLR*, 2020.
- [9] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo, “Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows,” in *Proceedings of ICCV*, 2021, pp. 10012–10022.
- [10] Ivan Sudakow, Vijayan K Asari, Ruixu Liu, and Denis Demchev, “MeltPondNet: A Swin Transformer U-Net for Detection of Melt Ponds on Arctic Sea Ice,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 8776–8784, 2022.
- [11] Junyuan Yao and Shuanggen Jin, “Multi-Category Segmentation of Sentinel-2 Images Based on the Swin UNet Method,” *Remote Sensing*, vol. 14, no. 14, pp. 3382, 2022.
- [12] Hongwei Dong, Lamei Zhang, and Bin Zou, “Exploring Vision Transformers for Polarimetric SAR Image Classification,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–15, 2021.
- [13] Anthony Fuller, Koreen Millard, and James R Green, “SatViT: Pretraining Transformers for Earth Observation,” *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2022.
- [14] Yuan Yuan and Lei Lin, “Self-Supervised Pretraining of Transformers for Satellite Image Time Series Classification,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 14, pp. 474–487, 2020.
- [15] Jørgen Buus-Hinkler, Tore Wulf, Andreas Rønne Stokholm, Anton Korosov, Roberto Saldo, and Leif Toudal Pedersen, “AI4Arctic Sea Ice Challenge Dataset,” *Technical University of Denmark*.
- [16] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang, “CVT: Introducing Convolutions to Vision Transformers,” in *Proceedings of ICCV*, 2021, pp. 22–31.
- [17] Sergey Ioffe and Christian Szegedy, “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift,” in *Proceedings of ICML*. PMLR, 2015, pp. 448–456.
- [18] Vinod Nair and Geoffrey E Hinton, “Rectified Linear Units Improve Restricted Boltzmann Machines,” in *Proceedings of ICML*, 2010.
- [19] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár, “Focal Loss for Dense Object Detection,” in *Proceedings of ICCV*, 2017, pp. 2980–2988.
- [20] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, “U-net: Convolutional Networks for Biomedical Image Segmentation,” in *Proceedings of MICCAI*. Springer, 2015, pp. 234–241.