



\*Correspondence:  
Mirakram Aghalarov, Baku  
Higher Oil School, Baku,  
Azerbaijan, [info@bhos.edu.az](mailto:info@bhos.edu.az)

# Utilization of Temporal Dimension in Satellite Imagery: Better Semantic Segmentation With Low Data Resources

Mirakram Aghalarov

*Baku Higher Oil School, Baku, Azerbaijan, [info@bhos.edu.az](mailto:info@bhos.edu.az)*

## Abstract

Time series image processing, a subfield of computer vision, enhances the accuracy of applications by leveraging temporal context. While this advantage is commonly utilized in video-based tasks, satellite imagery can also be treated as time series data when geospatial coordinates and timestamps are considered. Semantic segmentation, a key task in remote sensing, can benefit significantly from this temporal information. However, acquiring high-quality labeled datasets for such tasks remains a major challenge. In this study, we propose a novel temporal-aware domain adaptation framework for semantic segmentation, specifically targeting the detection of oil spills in the Caspian Sea. Our approach integrates time series information to improve cross-domain generalization. We evaluate our method on the synthetic SynthOil dataset, and a custom-labeled real-world dataset provided by Azercosmos and ArcGIS. Furthermore, we enhance the backbone of the Segformer model using a super-resolution dataset curated from Azercosmos and open data from the Esri ArcGIS platform. Experimental results demonstrate the effectiveness of our approach in improving segmentation performance across domains.

**Keyword:** Semantic Segmentation, Satellite Imagery, Deep Learning, Computer Vision, Temporal Dimension, Spatio-temporal Processing.

## 1. Introduction

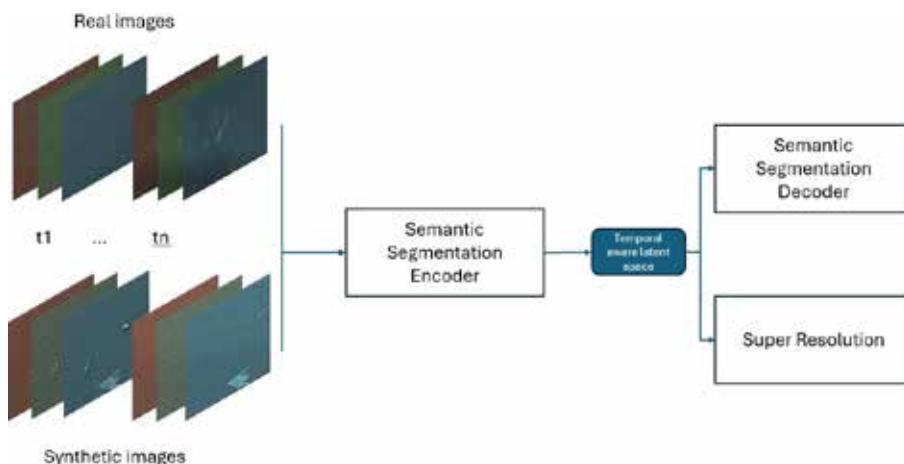
Video processing is challenging task in Computer Vision which constrains the application with their computational requirements. Considering given difficulty, temporal dimension in video processing can also bring several advantages such as more accurate results. These results are concluded in paper (Zhou, T., Porikli, F., Crandall, D. J., Van Gool, L., & Wang, W., 2022). which shows that using Semantic Segmentation with Spatio-Temporal features can increase the speed of the processing and accuracy at the same time.

On the other hand, Semantic Segmentation has wide range of applications from autonomous driving (Yao, S., Guan, R., et al., 2023) to land cover and land use (Ajibola, S., & Cabral, P., 2024). Considering the nature of applications, there is need

for efficient computations because of edge devices. Utilizing temporal dimensions in video processing increases efficiency in those applications. However, use cases like land cover- land use don't have any computational boundaries as processes are not executed in real-time. Therefore, video processing techniques can be used for increasing the accuracy of the model.

Accuracy of Semantic Segmentation models heavily relies on amount and cleanliness of the dataset which requires huge time consumption for labelling. Given downstream task requires accurate pixel annotations according to boundaries of the objects. This process can take up to 90 minutes per image (Schwonberg, M., Niemeijer, J., & Termöhlen, J. A., 2023). Another problem with the dataset is that it is difficult to cover all cases like anomalies or making them balanced. For example, our research objective, which is Oil Spill detection on Sea environments, has this certain problem. The reason is that there can be 2 or 3 cases in every 100 images which contain oil spill anomaly on the water. It also leads to additional time consumption for the search of anomaly cases, and it is not guaranteed that case will be found. These problems slow down the development of the custom applications even just with fine-tuning.

Solution to dataset problem is to use custom designed synthetic dataset for semantic segmentation (Schwonberg, M., Niemeijer, J., & Termöhlen, J. A., 2023). Every case and classes should be designed carefully so that they can represent the real world. The main reason for the synthetic dataset is automated annotation. Considering that in 3D environment creation software we can assign the object classes, corresponding ground truth logits can be generated with those software applications. However, it brings domain shift problems between the real and synthetic environment. Therefore, domain shift is solved by do- main adaptation methodologies so that the semantic segmentation model is trained on synthetic dataset and real dataset together but with only the label of the synthetic dataset.



*Fig. 1. Training pipeline for the semantic Segmentation model with Domain Adaptation. Latent space will be generated according to the differences in temporal dimension.*

Our paper proposes novel unsupervised domain adaptation (Figure 1) pipeline based on satellite imagery. Our assumption is that super resolution and temporal dimension together would help to make the backbone of the semantic segmentation more robust even if the model is trained on mainly with synthetic dataset. To apply this pipeline, we utilize 3 datasets: SynthOil, Azercosmos and Esri Sentinel hub open satellite imagery. Our novel unsupervised domain adaptation method surpasses the state-of-the-art methods in accuracy in real domain when there is low number of unlabelled datasets.

The research study contains the following steps:

- Analysis of datasets - Synthetic dataset formulation, real dataset collection and annotation for the validation. Some of the real datasets had different additional band information which could be beneficial for model development
- Baseline training with and without state-of-the-art domain adaptation methods - To understand the added value of the proposed approach, we benchmarked SOTA methods.
- Testing proposed approach - We made ablation studies with different steps of our approach so that we were able to understand contribution of each part.

## ***2. Related word***

### ***2.1. Semantic Segmentation.***

Semantic segmentation is a pixel-wise classification task, with models generally categorized into Encoder-Decoder and Two-Path structures. Encoder-Decoder models progressively reduce feature map sizes before reconstructing the original image, while two-path architectures balance feature extraction and efficiency.

One of the earliest and most influential models, UNet (Ronneberger, O., Fischer, P., & Brox, T., 2015, October), introduced a U-shaped design with skip connections to enhance feature retention. However, its computational inefficiency led to the development of more optimized models. DeepLab (Chen, L. C., Papandreou, G., Schroff, F., & Adam, H., 2017) significantly improved performance by integrating Atrous Spatial Pyramid Pooling (ASPP), reducing memory usage while maintaining accuracy, an essential step for real-time applications.

For latency-sensitive tasks like autonomous driving (Yao, S., Guan, R., et al., 2023), models such as BiSeNet (Yu, C., Gao, C., Wang, J., Yu, G., Shen, C., & Sang, N., 2021), Fast-SCNN have been widely adopted (Poudel, R. P., Liwicki, S., & Cipolla, R., 2019). BiSeNet (Yu, C., Gao, C., Wang, J., Yu, G., Shen, C., & Sang, N., 2021) introduced a two-path structure to maintain a high receptive field while keeping parameter count low. Its “context path” captures global information efficiently, while the “spatial path” focuses on detailed feature extraction, achieving a balance between speed and accuracy. Fast-SCNN follow similar principles to optimize segmentation for real-time deployment (Poudel, R. P., Liwicki, S., & Cipolla, R., 2019).

With the rise of Vision Transformers (ViTs) (Dosovitskiy, A., 2020), transformer-based segmentation models have outperformed CNNs, particularly in domain-adaptive tasks.

ViTs require extensive training data, but advancements like DEiT (Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., & Jégou, H., 2021, July) (with token distillation) and LeViT (optimized for real-time inference) have improved their efficiency (Graham, B., El-Nouby, A., et al., 2021). Transformer-based architecture consistently shows superior mean Intersection over Union (mIoU) scores compared to CNNs. SegFormer (Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J. M., & Luo, P., 2021), a widely used transformer-based model, combines a lightweight MiT (Mixed Transformer) encoder with a CNN-based decoder. DaFormer builds on this by integrating ASPP into the same backbone, demonstrating further accurate improvements (Hoyer, L., Dai, D., & Van Gool, L., 2022). Given resource constraints, selecting a model suited for real-time applications while balancing performance and memory efficiency remains a critical consideration.

## *2.2. Domain Adaptation.*

Semantic segmentation models are highly sensitive to domain shifts, categorized into label shift, concept shift, conditional shift, and covariate shift (Schwonberg, M., Niemeijer, J., & Termöhlen, J. A., 2023). Our primary focus is on covariate and conditional shifts, which arise in dynamic environments. Label shift occurs when class distributions vary across datasets, while concept shifts stem from differences in grouping conventions across regions. Various studies have explored methods to address these challenges, with Unsupervised Domain Adaptation (UDA) being a common approach. UDA techniques primarily fall into feature-level and pixel-level adaptation. Feature-level adaptation improves robustness in feature extraction, while domain-level adaptation aligns features across datasets for better generalization.

Prototypical Learning - Few-shot learning methods have shown strong results, particularly on limited datasets. Cross-domain Prototypical Learning has been effective in refining domain spaces, as demonstrated by Bi-directional Contrastive Learning (Lee, G., Eom, C., Lee, W., Park, H., & Ham, B., 2022, October), which enhances intra-class compactness while increasing inter-class separation. This approach also addresses limitations of self-training by utilizing dynamic pseudo-labels to prevent overfitting.

Pixel-Level Adaptation - modifies input images before training, ensuring that synthetic and real domains share similar characteristics. AdvStyle applies augmentation techniques for domain adaptation, while Fourier Domain Adaptation transfers high-frequency and amplitude features between domains. ProCST enhances UDA performance by employing a cyclic style transfer mechanism, reducing domain discrepancies before training.

Self-training - has further advanced domain adaptation by refining pseudo-labeling techniques. Some methods introduce pre-training stages to stabilize pseudo-labels before applying contrastive learning, as seen in SimCLR-based approaches (Chen, T., Kornblith, S., Norouzi, M., & Hinton, G., 2020, November). DaFormer integrates transformers into segmentation models while incorporating strategies like rare-class

sampling and feature distance computation (Hoyer, L., Dai, D., & Van Gool, L., 2022). To mitigate pseudo-label reliability issues, an Exponential Moving Average update is applied to the teacher network.

Overall, these advancements collectively enhance the adaptability of segmentation models to domain shifts, contributing to more reliable and generalizable performance across diverse environments.

### **2.3. Super Resolution.**

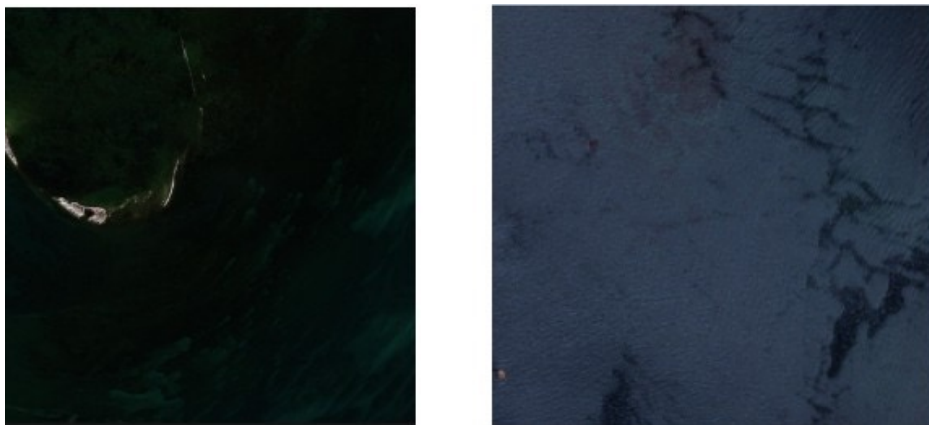
For making the encoder of the model more robust on real dataset we plan to use super-resolution header which should be compatible to our backbone mechanism. SwinIR (Liang, J., Cao, J., et al., 2021) proposes a novel architecture for image restoration with shallow and deep feature extraction as (Hsu, C. C., Lee, C. M., & Chou, Y. S., 2024). Considering that our backbone is based on attention mechanism (Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J. M., & Luo, P., 2021) with mixed transformers, focusing with (Liang, J., Cao, J., et al., 2021) would be wise for better compatibility. However, main purpose is not to have highly accurate super resolution pipeline. It is the tool for forcing the model backbone for training in real domain.

## **3. Methodology**

### **3.1. Dataset.**

Several considerations have been made for dataset collection. The main aspect which is going to affect accuracy is temporal dimensions. Time-series feature can be obtained by using GeoTIFF data from different timestamps. Therefore, we have collected data accordingly:

- Azercosmos - For our research, Azercosmos has provided 500 km<sup>2</sup> from Caspian Sea (Figure 2). We have been provided with 5 samples, from 5 years. Each sample contains 1m of resolution with 4 bands. Those bands - Red, Green, Blue, Near-Infrared

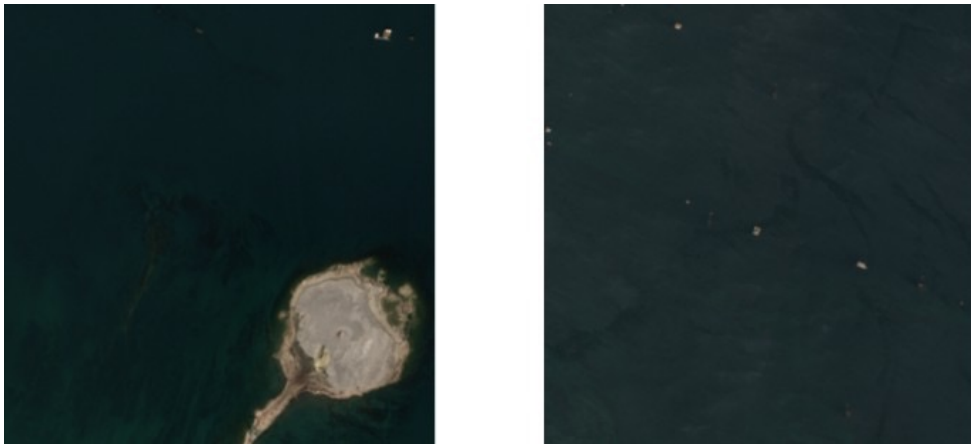


*Fig. 2. Azercosmos Dataset - Dataset consists of 5 years of captures from Caspian Sea with 1 m of resolution*

- have been analyzed well and we selected only RGB bands for the training. The main reason is that NIR signals do not have penetration in water which does not let us classify oil in our application. Additional limitations are dependence on sunlight and confusion with other substances. However, we can utilize near-infrared signals for further post processing to identify oil sickness over the water. Dataset will be used for super resolution and pixel - level domain adaptation (Ettegui, S., Abu-Hussein, S., & Giryes, R., 2022).

- Sentinel Data - Esri and ArcGIS platform provide easy to use large-scale, updated, open-source satellite imagery for research purposes. For successful super resolution application, we focused on getting the same areas with the Azercosmos as Sentinel Data is 10 m resolution (Figure 3). However, we can access satellite imagery every month for 10 years. Each image contains 13 bands for different purposes like vegetation, moisture, and emission control. Dataset will be used for Pixel-level domain adaptation

- ProCST training (Ettegui, S., Abu-Hussein, S., & Giryes, R., 2022)

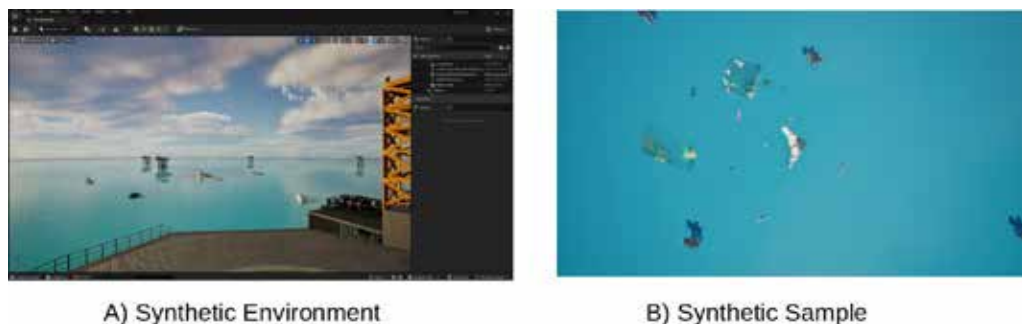


**Fig. 3. Sentinel Dataset collected from Esri - Dataset consist of 10 years of captures from Caspian Sea with 10 m of resolution**

- SynthOil - is a synthetic dataset created with Unreal Engine. Dataset contains 15 scenes in which several dynamic factors are considered: Sun angle, ship position, leak position. As every scene has multiple images from the same position, we can consider them as time series data. Therefore, even in our domain adaptation pipeline we can utilize this situation. Unreal Engine did not have sufficient tools for generating NIR or other band information. Therefore, on all our datasets we will be focusing on RGB channels. Dataset contains the instance segmentation labels. 4 classes are available: Ship - small and dynamic objects, Oil Platform - Less dynamic, mostly static objects, Oil Spills - Main objective and dynamic object, Soil - Big classes and static objects. As there are differences in shapes, there is definite class imbalance. Example of the

dataset is shown in Figure 4.

- Validation dataset - for validation of the different approaches we labelled 200 real satellite image tiles for Semantic Segmentation. Dataset consists of 100 images from Azercosmos, 100 images from Sentinel satellites.



*Fig. 4. SynthOil Dataset based on Unreal Engine 5. A) Unreal Engine development environment. B) Satellite view for unreal engine environment*

### 3.2. Training Procedure.

Our methodology requires dividing the model into backbone and header parts for better analysis in latent space. We must analyze the output of each block in different layers for each Synthetic sample and real samples:

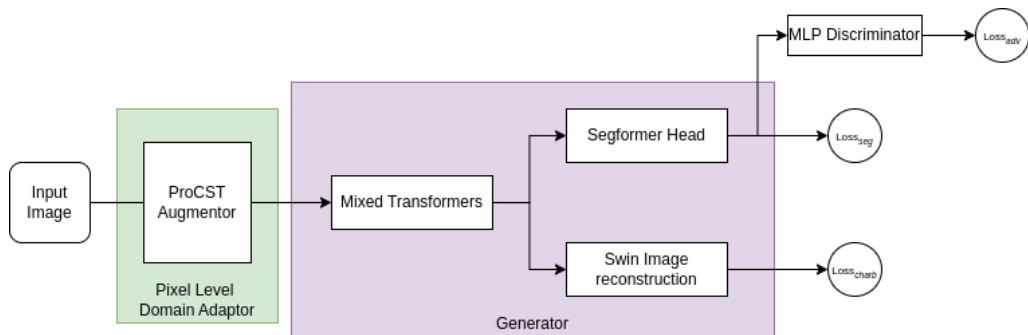
- Synthetic and real images - To decrease the domain shift in the pixel space, rule-based and deep learning-based transformations should be applied to images.
- Temporal Dimension - Data should contain the information from the last capture (previous month or year does not matter) so that changes between the captures should generate more sophisticated features.
- Extracted Features from backbone - latent values of input from synthetic and real images extracted by the backbone should be similar for robustness of the decoder.
- Semantic Segmentation prediction - Adversarial attack should be applied for better output to control semantic segmentation.

By controlling domain shift in each layer mentioned above, better domain adaptation method and higher accuracy can be achieved. Our pipeline is composed of the following items and demonstrated in Figure 5:

- ProCST module - This deep transformation should allow us to close the gap between the real and synthetic images.
- Semantic Segmentation backbone - Mixed Transformers as given in Segformer (Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J. M., & Luo, P., 2021). This backbone will remain the same for the semantic Segmentation. This backbone (we can call encoder too) is most important part of the model as extracted features are key points for the preliminary robustness.



- Segformer decoder - Downstream task specific module which will generate desired se- mantic maps. This module consists of MLP layers which decrease the number of heavy computations for generation of the semantic map
- SwinIR decoder - Super resolution part which is essential for adapting the encoder for forcing the model to generate essential features for the real domain too as super resolution dataset is from real environment.
- Discriminator - it is the main module for adversarial attacks. Main purpose of this block is that semantic segmentation prediction will be checked if the output generated from synthetic or real.



*Fig. 5. Architecture of Domain Adaptation methodology: Composed of Do- main adaptor, Generator and discriminator.*

As first stage of training, pixel-level domain adaptation is applied to change the synthetic images to make them more alike real images. In other words, patterns of real images are transferred into synthetic images. It is done through utilization of ProCST which contains converting synthetic images to real and real to synthetic domains (Etteedgui, S., Abu-Hussein, S., & Giryes, R., 2022). It allows us to learn patterns for both domains and get more accurate results.

After getting results from ProCST architecture, proposed temporal dimensions can be utilized in input. Previous captures (also transformed by ProCST) are concatenated as channels and fed into the segformer architecture. In every iteration, SynthOil batch is fed and Dice loss (Equation 1) calculated for the label of the synthetic dataset.  $p_i$  is the predicted probability for pixel  $i$ ,  $g_i$  is the ground truth label (binary) for pixel  $i$ ,  $N$  is the total number of pixels in the image is a small constant to prevent division by zero.

$$\mathcal{L}_{Dice} = 1 - \frac{2 \sum_{i=1}^N p_i g_i}{\sum_{i=1}^N p_i + \sum_{i=1}^N g_i + \epsilon} \quad (1)$$

2nd part of the iteration is based on super resolution, in which dataloader is combining Synthoil and Azercosmos dataset by decreasing resolution to create input-label relationship. Loss for the super resolution is the Charbonnier Loss (Equation 2) from the paper (Liang, J., Cao, J., et al., 2021) which is defined as:



$$\mathcal{L}_{charb} = \frac{1}{N} \sum_{i=1}^N \sqrt{(I_{HR}^i - I_{SR}^i)^2 + \epsilon^2} \quad (2)$$

where:

$I_{HR}^i$  is the ground truth high-resolution pixel value,

$I_{SR}^i$  is the super-resolved output pixel value,

- $\epsilon$  is a small constant ( $10^{-3}$ ) to prevent numerical instability,
- $N$  is the total number of pixels in the image.

The last part of the training procedure involves the adversarial part in which Sentinel data is fed to the model and the output of the semantic segmentation model is checked by discriminator if it comes from real or synthetic image. Whether model is able to deceive the discriminator or not. This is done through adversarial loss (Equation 3).

$$\mathcal{L}_{adv} = \mathbb{E}_{x_s \sim P_s} [\log D(x_s)] + \mathbb{E}_{x_t \sim P_t} [\log (1 - D(x_t))] \quad (3)$$

where:

- $x_s$  is a sample from the source domain distribution  $P_s$ ,
- $x_t$  is a sample from the target domain distribution  $P_t$ ,
- $D(x)$  is a domain discriminator that predicts whether  $x$  is from the source or target domain.

Several stages of the experiments have been carried out to understand the effectiveness of the method. State-of-the-art unsupervised domain adaptation methods were selected as a for comparison. Additionally, CNN based architecture is also tested to understand the effect of the proposed pipeline.

As a state-of-the-art method we selected Fourier Domain Adaptation (Yang, Y., & Soatto, S., 2020), Daformer, Bidirectional Contrastive Learning to train with our selected datasets (Lee, G., Eom, C., Lee, W., Park, H., & Ham, B., 2022, October). All experiments have been done in distributed mode with 4x Nvidia RTX4090 graphics card. Results are demonstrated in Table 1.

As it is seen from Table 1, proposed method performs the best in given comparison with the state-of-the-art methods in mIoU measurements, Oil Spills and Oil platform classes. Baseline shows that there is huge domain shift between the synthetic and

Model Name	Ship	Oil Plateorm	Oil Spill	Soil	mIoU
Baseline	0.0	28.1	18.2	66.1	28.1
ProCST	5.8	31.0	32.2	75.4	36.1
DaFormer	41.1	48.3	68.7	89.1	61.8
Bidirectional Contrastive L.	46.8	51.2	62.0	74.4	58.4
Fourier Domain Adaptation	36.1	40.0	51.2	67.5	48.7
Ours (with Segeormer)	42.1	52.1	71.3	86.9	63.1

**Table 1. Comparison of the proposed method with the state-of-the-art methods with intersection over union score (IoU) per class and Mean Intersection over Union.**

Model Name	ProCST	Super resolution	Adversarial	Temporal Dimension	mIoU
Segeormer B3					28.1
Segeormer B3	+				36.1
Segeormer B3	+	+			48.2
Segeormer B3	+	+	+		60.0
Segeormer B3	+	+	+	+	63.1
DeepLab v3					19.9
DeepLab v3	+				23.2
DeepLab v3	+	+			41.5
DeepLab v3	+	+	+		48.9
DeepLab v3	+	+	+	+	55.0

**Table 2. Ablation study based on 2 different models: results are compared according to the mean Intersection over Union in %**

real domains. Bidirectional Contrastive learning is surpassing the performance of ours in small classes like ship as this training method relies on cropping, augmenting and comparison. Very big objects like soil are classified better by Daformer because of self-supervised nature.

To understand the behavior of each component separately, we have carried out ablation study which is demonstrated in Table 2.

As is seen from Table 2, every component contributes to the result positively. For measurement of effectiveness, we have also used Deeplab (Chen, L. C., Papandreou, G., Schroff, F., & Adam, H., 2017) architecture with the backbone of ResNet50 (He, K., Zhang, X., Ren, S., & Sun, J., 2016). It is clear from the results that convolution-based architecture is not robust to domain changes. Self-Attention-based architecture surpasses CNN with high margin.

## 5. Conclion

The proposed approach proves its effectiveness for given application with satellite imagery. With this technique, it is possible to get temporal-aware latent space based on time-series syn- thetic and real dataset. Considering impact of the application, other critical use cases can also be applied with minimal time consumption. According to the results, high accuracy of the transformers models contributed to the success of our methodology.

## 6. Limitations

Despite the advantages of the methodology, there are several limitations which should be considered:

- Pretrained models are not applicable to our pipeline because of the input changes. This can increase the computational cost as pretraining is computationally expensive.
- Training data should be carefully collected and needs to be time series. For super resolution, high quality data is necessary which decreases the possibilities to find.

### Reference

Ajibola, S., & Cabral, P. (2024). A systematic literature review and bibliometric analysis of semantic segmentation models in land cover mapping. *Remote Sensing*, 16(12), 2222.

Chen, L. C., Papandreou, G., Schroff, F., & Adam, H. (2017). Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*.

Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020, November). A simple framework for contrastive learning of visual representations. *In International conference on machine learning* (pp. 1597-1607). PmLR.

Dosovitskiy, A. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.

Ettehadgui, S., Abu-Hussein, S., & Giryas, R. (2022). ProCST: Boosting semantic segmentation using progressive cyclic style-transfer. *arXiv preprint arXiv:2204.11891*.

Graham, B., El-Nouby, A., et al. (2021). Levit: a vision transformer in convnet's clothing for faster inference. *In Proceedings of the IEEE/CVF international conference on computer vision* (pp. 12259-12269).

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *In Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).

Hoyer, L., Dai, D., & Van Gool, L. (2022). Daformer: Improving network architectures and training strategies for domain-adaptive semantic segmentation. *In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 9924-9935).

Hsu, C. C., Lee, C. M., & Chou, Y. S. (2024). Drct: Saving image super-resolution away from information bottleneck. *In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 6133-6142).

Lee, G., Eom, C., Lee, W., Park, H., & Ham, B. (2022, October). Bi-directional contrastive learning for domain adaptive semantic segmentation. *In European Conference on Computer Vision* (pp. 38-55). Cham: Springer Nature Switzerland.

Liang, J., Cao, J., et al. (2021). Swinir: Image restoration using swin transformer. *In Proceedings of the IEEE/CVF international conference on computer vision* (pp. 1833-1844).

Poudel, R. P., Liwicki, S., & Cipolla, R. (2019). Fast-scnn: Fast semantic segmentation network. *arXiv preprint arXiv:1902.04502*.

Ronneberger, O., Fischer, P., & Brox, T. (2015, October). U-net: Convolutional networks for biomedical image segmentation. *In International Conference on Med-*

*ical image computing and computer-assisted intervention* (pp. 234-241). Cham: Springer international publishing.

Schwonberg, M., Niemeijer, J., & Termöhlen, J. A. (2023). Survey on unsupervised domain adaptation for semantic segmentation for visual perception in automated driving. *IEEE Access*, 11(54296-54336), 1-2.

Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., & Jégou, H. (2021, July). Training data-efficient image transformers & distillation through attention. *In International conference on machine learning* (pp. 10347-10357). PMLR.

Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J. M., & Luo, P. (2021). SegFormer: Simple and efficient design for semantic segmentation with transformers. *Advances in neural information processing systems*, 34, 12077-12090.

Yang, Y., & Soatto, S. (2020). Fda: Fourier domain adaptation for semantic segmentation. *In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 4085-4095).

Yao, S., Guan, R., et al. (2023). Radar-camera fusion for object detection and semantic segmentation in autonomous driving: A comprehensive review. *IEEE Transactions on Intelligent Vehicles*, 9(1), 2094-2128.

Yu, C., Gao, C., Wang, J., Yu, G., Shen, C., & Sang, N. (2021). Bisenet v2: Bilateral network with guided aggregation for real-time semantic segmentation. *International journal of computer vision*, 129(11), 3051-3068.

Zhou, T., Porikli, F., Crandall, D. J., Van Gool, L., & Wang, W. (2022). A survey on deep learning technique for video segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 45(6), 7099-7122.

**Submitted: 19.05.2025**

**Accepted: 08.09.2025**