

Towards Explainable Multimodal Land Cover Segmentation Using Swin Transformer

Islomjon Shukhratov^{*1}, Mehak Khan¹, and Reza Arghandeh¹

¹Western Norway University of Applied Science

{islomjon.shukhratov, mehak.khan, reza.arghandeh}@hvl.no

Abstract

Recent advancements in Vision Transformers (ViTs) demonstrate strong potential in remote sensing, providing powerful spatial feature representations for complex land cover segmentation tasks. In this study, we explore multimodal data fusion of Synthetic Aperture Radar (SAR) and optical imagery for land cover mapping. We train and evaluate Swin Transformer models and employ explainable AI (xAI) techniques to analyse the contribution of each modality and feature to the model's predictions. We expect to improve the interpretability and robustness of multimodal remote sensing models for land cover segmentation.

1 Introduction

Recent advancements in ViTs show state-of-the-art performance in computer vision tasks such as classification and segmentation by leveraging self-attention mechanisms to learn strong representations of spatial features in images [1]. In remote sensing applications, complex land cover patterns and diverse terrain characteristics pose significant challenges, ViTs demonstrate a promise in addressing these challenges, particularly in building segmentation, change detection, scene classification [2, 3].

Multimodal data fusion refers to a combination of several types of data in order to enhance the performance of the deep learning algorithms. In land cover segmentation, the most popular sources for multimodal data fusion are optical and SAR images [4]. However, only a limited number of studies explore the use of transformers for SAR-optical data fusion in land cover classification tasks. In this project, we aim to address this gap by using Sentinel-1 and Sentinel-2 data to train models in one geographic region (the Grand-Est region of France) and evaluate their generalisation in another (Askvoll municipality, Norway). Additionally, we plan to investigate the explainability of the proposed multimodal fusion approach to better understand the model's decision making process.

^{*}Corresponding Author.

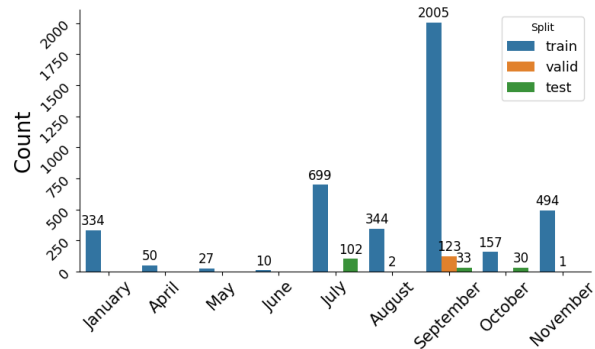


Figure 1. Number of samples by month.

2 Dataset

We use MultiSenGE dataset provided by Wenger et al. [5]. The dataset covers a large territory in east of the France with the area of $57,433 \text{ km}^2$ and consists of Sentinel-1 and Sentinel-2 images for 2020 year, and regional land use and land cover maps provided by OCSGE2-GEOGRANDEST for 2019/2020 years. In total, there are 8157 unique patches of 256×256 pixels generated from 14 tiles of Sentinel-2, and the dataset has 1,012,22 patches for Sentinel-1 and 72,033 patches for Sentinel-2. There are 14 classes in the dataset, but we remap them into 4 main categories to avoid data imbalance issue: urban, agricultural, forest and water areas.

To facilitate the computation resource constraints, we select only the patches that have matching capture date for both sensors, and the resulting number of images is 4411. We split dataset into train, validation and test sets based on tiles, rather than a random split to avoid data leakage and contamination. Consequently, there are 4121, 127 and 167

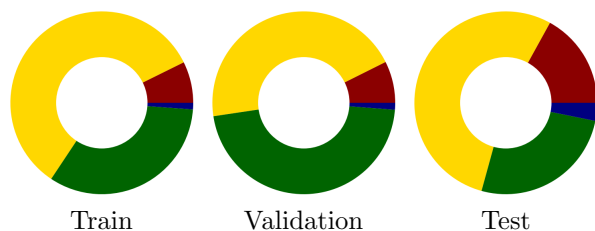


Figure 2. Class distribution, red: urban, yellow: agriculture, green: forest, blue: water areas.

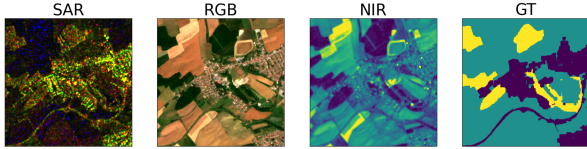


Figure 3. Example of the images.

063 patches for train, validation and test sets. Figure 1
 064 shows monthly distribution of the selected patches.
 065 Additionally, we plan to evaluate the generaliza-
 066 tion capability of the transformer models on unseen
 067 data and different geographic conditions. We des-
 068 ignate Askvoll municipality in the Western part of
 069 Norway which mainly consists of forests and small
 070 lakes in the mountainous terrain. Similar to Multi-
 071 SenGE dataset, we collect Sentinel-1 and Sentinel-2
 072 images for the 2022 year, and we use the land cover
 073 map provided by Norwegian Institute of Bioeconomy
 074 (NIBIO) as a ground truth data.

075 3 Methodology

076 In our experiments, we intent to use the Shifted
 077 Window (Swin) Transformer. Swin Transformer
 078 is a hierarchical ViT architecture that computes
 079 self-attention within shifted local windows, enabling
 080 efficient modelling of both local and global visual de-
 081 pendencies while maintaining linear computational
 082 complexity with image size [6]. SWIN transformer
 083 is available at small, medium and large sizes, and we
 084 train all models for better understanding their capa-
 085 bility. Furthermore, to benchmark our results, we
 086 aim to train models based on Convolutional Neural
 087 Networks (CNN), namely U-Net [7] and DeepLabV3
 088 [8] architectures.

089 We plan to utilise the VV and VH polarisation
 090 bands of the SAR data along with the RGB and
 091 NIR bands of the optical imagery, and Normalized
 092 Difference Vegetation Index (NDVI) as an additional
 093 input feature for training models to perform land
 094 cover classification. Figure 3 depicts a sample from
 095 the training dataset.

096 We also plan to investigate the explainability of
 097 the SAR–optical multimodal fusion using xAI tech-
 098 niques. Our goal is to understand how each modality
 099 and feature contributes to the model’s land cover
 100 classification decisions, providing insights into the
 101 relative importance of SAR and optical inputs and
 102 improving the interpretability of multimodal remote
 103 sensing models. Figure 4 displays the pipeline of
 104 proposed model.

105 4 Expected Results

106 To evaluate the models’ performance, we compute
 107 the F1-score and Intersection over Union (IoU) score

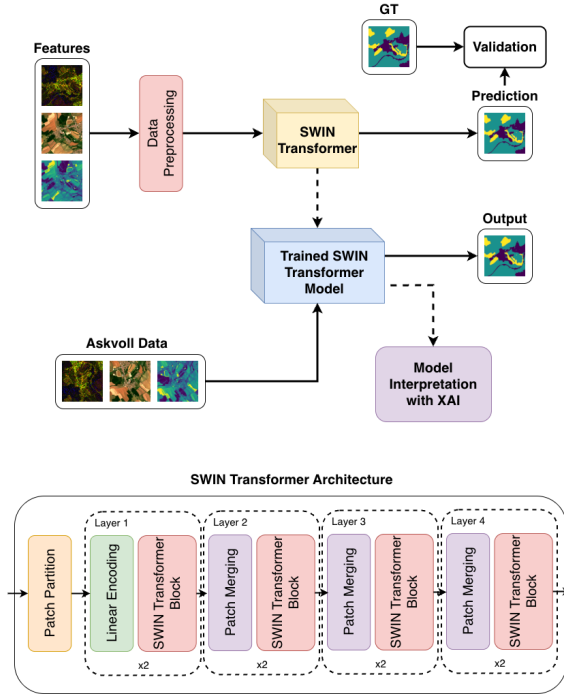


Figure 4. Pipeline of the proposed model.

between the predicted image and the ground truth
 108 map computed by the following formula: 109

$$110 \quad F1 \text{ score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (1)$$

$$111 \quad \text{IoU} = \frac{|\text{Prediction} \cap \text{Ground Truth}|}{|\text{Prediction} \cup \text{Ground Truth}|} \quad (2)$$

112 We define acceptable performance as a class-
 113 average IoU above 70% and F1-score above 80%, We
 114 anticipate notable challenges during domain shift
 115 evaluation, as the contrasting geographic and
 116 environmental conditions. Therefore, we expect 10%
 117 drop in metrics compared to the trained dataset.

118 For xAI interpretability, we apply methods such
 119 as Grad-CAM [9], Captum [10] and attention map
 120 visualization to assess the spatial relevance and
 121 contribution of SAR and optical features to the model’s
 122 decisions. This is important to understand which
 123 modality dominates in a specific class prediction.

124 5 Conclusion

125 We present an ongoing study on multimodal land
 126 cover classification using SAR and optical imagery
 127 with vision transformers. The proposed approach
 128 demonstrates the potential of multimodal fusion to
 129 improve segmentation performance, and future ex-
 130 plainability analyses using xAI will provide deeper
 131 insights into the contribution of each modality. These
 132 results lay the groundwork for developing more in-
 133 terpretable and robust remote sensing models for
 134 land cover mapping.

135 **References**

- 136 [1] K. Han, Y. Wang, H. Chen, X. Chen, J. Guo,
137 Z. Liu, Y. Tang, A. Xiao, C. Xu, Y. Xu, Z.
138 Yang, Y. Zhang, and D. Tao. “A Survey on
139 Vision Transformer”. In: *IEEE Transactions*
140 *on Pattern Analysis and Machine Intelligence*
141 45.1 (2023), pp. 87–110. DOI: [10.1109/TPAMI.](https://doi.org/10.1109/TPAMI.2022.3152247)
142 [2022.3152247](https://doi.org/10.1109/TPAMI.2022.3152247).
- 143 [2] Z. Xu, W. Zhang, T. Zhang, Z. Yang, and
144 J. Li. “Efficient transformer for remote sens-
145 ing image segmentation”. In: *Remote Sensing*
146 13.18 (2021), p. 3585.
- 147 [3] S. Hao, B. Wu, K. Zhao, Y. Ye, and W. Wang.
148 “Two-stream swin transformer with differen-
149 tiable sobel operator for remote sensing image
150 classification”. In: *Remote Sensing* 14.6 (2022),
151 p. 1507.
- 152 [4] J. Li, D. Hong, L. Gao, J. Yao, K. Zheng, B.
153 Zhang, and J. Chanussot. “Deep learning in
154 multimodal remote sensing data fusion: A com-
155 prehensive review”. In: *International Journal*
156 *of Applied Earth Observation and Geoinfor-*
157 *mation* 112 (2022), p. 102926. ISSN: 1569-8432.
158 DOI: [https://doi.org/10.1016/j.jag.](https://doi.org/10.1016/j.jag.2022.102926)
159 [2022.102926](https://doi.org/10.1016/j.jag.2022.102926).
- 160 [5] R. Wenger, A. Puissant, J. Weber, L.
161 Idoumghar, and G. Forestier. “Multisenge :
162 A Multimodal And Multitemporal Benchmark
163 Dataset For Land Use/land Cover Remote
164 Sensing Applications”. In: *ISPRS Annals of*
165 *the Photogrammetry, Remote Sensing and Spa-*
166 *tial Information Sciences V-3-2022* (2022),
167 pp. 635–640. DOI: [10.5194/isprs-annals-V-](https://doi.org/10.5194/isprs-annals-V-3-2022-635-2022)
168 [3-2022-635-2022](https://doi.org/10.5194/isprs-annals-V-3-2022-635-2022).
- 169 [6] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z.
170 Zhang, S. Lin, and B. Guo. “Swin transformer:
171 Hierarchical vision transformer using shifted
172 windows”. In: *Proceedings of the IEEE/CVF*
173 *international conference on computer vision.*
174 2021, pp. 10012–10022.
- 175 [7] O. Ronneberger, P. Fischer, and T. Brox.
176 “U-net: Convolutional networks for biomed-
177 ical image segmentation”. In: *International*
178 *Conference on Medical image computing and*
179 *computer-assisted intervention*. Springer. 2015,
180 pp. 234–241.
- 181 [8] L.-C. Chen, G. Papandreou, I. Kokkinos, K.
182 Murphy, and A. L. Yuille. “Deepplab: Sema-
183 ntic image segmentation with deep convolu-
184 tional nets, atrous convolution, and fully con-
185 nected crfs”. In: *IEEE transactions on pattern*
186 *analysis and machine intelligence* 40.4 (2017),
187 pp. 834–848.
- [9] R. R. Selvaraju, M. Cogswell, A. Das, R. 188
Vedantam, D. Parikh, and D. Batra. “Grad- 189
CAM: Visual Explanations from Deep Net- 190
works via Gradient-Based Localization”. In: 191
2017 IEEE International Conference on Com- 192
puter Vision (ICCV). 2017, pp. 618–626. DOI: 193
[10.1109/ICCV.2017.74](https://doi.org/10.1109/ICCV.2017.74). 194
- [10] N. Kokhlikyan, V. Miglani, M. Martin, E. 195
Wang, B. Alsallakh, J. Reynolds, A. Mel- 196
nikov, N. Kliushkina, C. Araya, S. Yan, and 197
O. Reblitz-Richardson. “Captum: A unified 198
and generic model interpretability library for 199
PyTorch”. In: *CoRR* abs/2009.07896 (2020). 200
arXiv: [2009.07896](https://arxiv.org/abs/2009.07896). URL: [https://arxiv.](https://arxiv.org/abs/2009.07896) 201
[org/abs/2009.07896](https://arxiv.org/abs/2009.07896). 202