

Vision Transformers

Konstantin A. Maslov
k.a.maslov@utwente.nl

University of Twente
Faculty of Geo-Information Science and Earth Observation

29 Oct 2022



Outline

Essential basics

- Layer normalization

- Multi-head self attention

ViT

Other architectures

- DeiT

- DPT

- SETR

- Segmenter

- Swin transformer

- SegFormer

- MLP-Mixer

Common practices

Summary

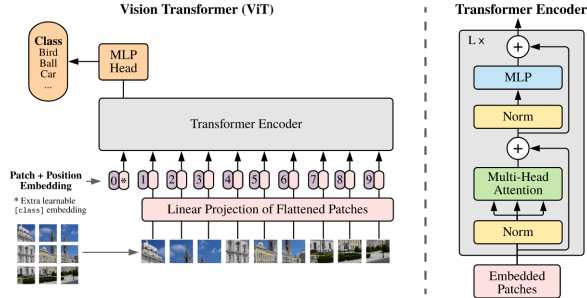


Layer normalization



Multi-head self attention

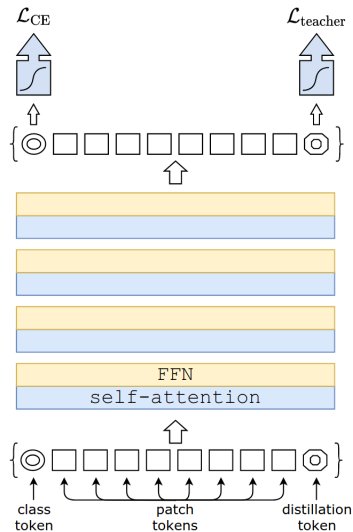




- ▶ Transformer design from NLP with minimal changes
- ▶ Achieves performance close to the state-of-the-art (CNNs) in classification tasks
- ▶ Can learn long-range relations in the very first layers due to the multi-head self-attention (which has quadratic time complexity)
- ▶ Requires pre-training on huge datasets to achieve good performance

<https://github.com/konstantin-a-maslov/transformers-seminar>

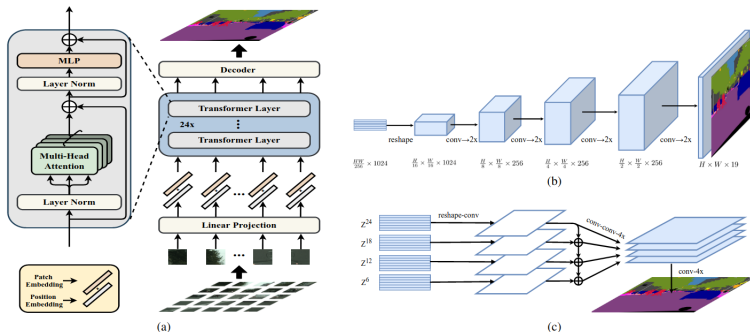




- ▶ ViT trained with distillation
- ▶ The teacher model is a CNN
- ▶ The authors claim that it reduces the amount of data required to train a transformer

Ranftl, R., Bochkovskiy, A., & Koltun, V. (2021). Vision Transformers for Dense Prediction. <https://doi.org/10.48550/arXiv.2103.13413>



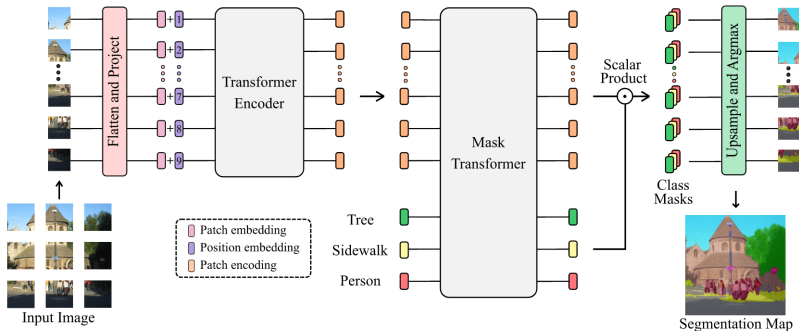


- ▶ ViT with a typical upsampling decoder as in FCNs
- ▶ The model has shown state-of-the-art performance in some tasks

Zheng, S., Lu, J., Zhao, H., Zhu, X., Luo, Z., Wang, Y., Fu, Y., Feng, J., Xiang, T., Torr, P. H. S., & Zhang, L. (2020). Rethinking Semantic Segmentation from a Sequence-to-Sequence Perspective with Transformers. Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 6877–6886. <https://doi.org/10.48550/arxiv.2012.15840>

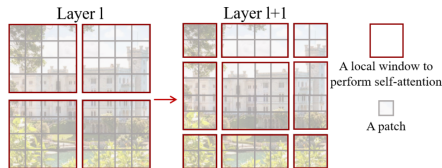
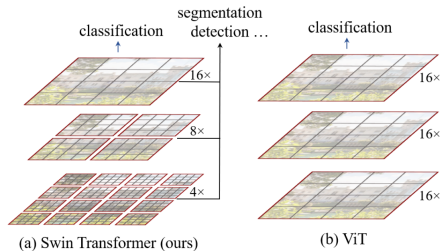


Segmenter



- ▶ Transformer-based decoder
- ▶ Introduced class embeddings
- ▶ The authors emphasized that transformer-based models are not so good at generating sharp object boundaries
- ▶ Does not seem to be a popular choice nowadays

Swin transformer

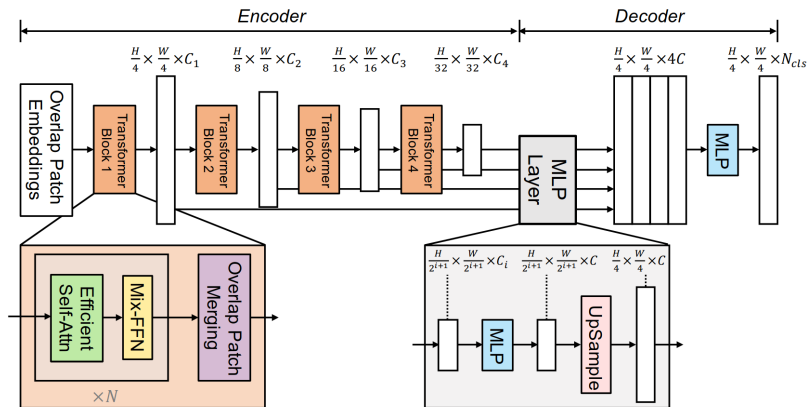


- ▶ Has linear time complexity due to the hierarchical design
- ▶ Introduced shifting windows

Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., & Guo, B. (2021). Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. Proceedings of the IEEE International Conference on Computer Vision, 9992–10002. <https://doi.org/10.48550/arxiv.2103.14030>

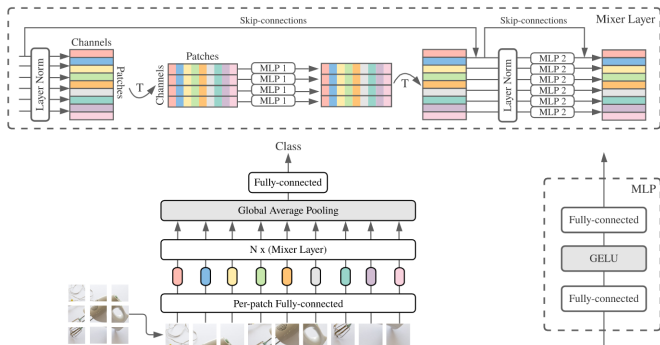


SegFormer



- ▶ The authors focused on an efficient design
- ▶ Seems to be the state-of-the-art among transformers for semantic image segmentation (general tasks) nowadays

MLP-Mixer



- Not a transformer!
- Replaces multi-head self-attentions with simple MLPs
- The authors have shown that it still possible to obtain good results with this design

Common practices

- ▶ 1
- ▶ 2
- ▶ 3



Summary

- ▶ Transformers require a lot of data to train
 - ▶ Which can be not our case
 - ▶ Pre-training is complicated due to the absence of huge datasets for multispectral/SAR/DEM data
 - ▶ Using the weights from more common datasets (ImageNet, JFT, ...) is still an option though, but it requires studies on how to better 'generalise' them for non-RGB images
- ▶ Seems like transformers are bad at restoring sharp boundaries in segmentation maps
 - ▶ Can be crucial as the spatial resolution of the imagery we use is very different from the resolution of the images from ImageNet or similar datasets
 - ▶ Smaller patch sizes improve the situation (if one has enough memory...)
 - ▶ Perhaps, there is a space to explore hybrid CNN-Transformer models
- ▶ However, as it will be shown later, a simple SETR-like transformer can still generate interesting results

