

# Foundation Models Permit Retinal Layer Segmentation Across OCT Devices

Olivier Morelle<sup>1,2</sup> and Thomas Schultz<sup>1,3</sup>

<sup>1</sup> University of Bonn, B-IT and Department of Computer Science

<sup>2</sup> Department of Ophthalmology, University Hospital Bonn, Bonn, Germany

<sup>3</sup> Lamarr Institute for Machine Learning and Artificial Intelligence

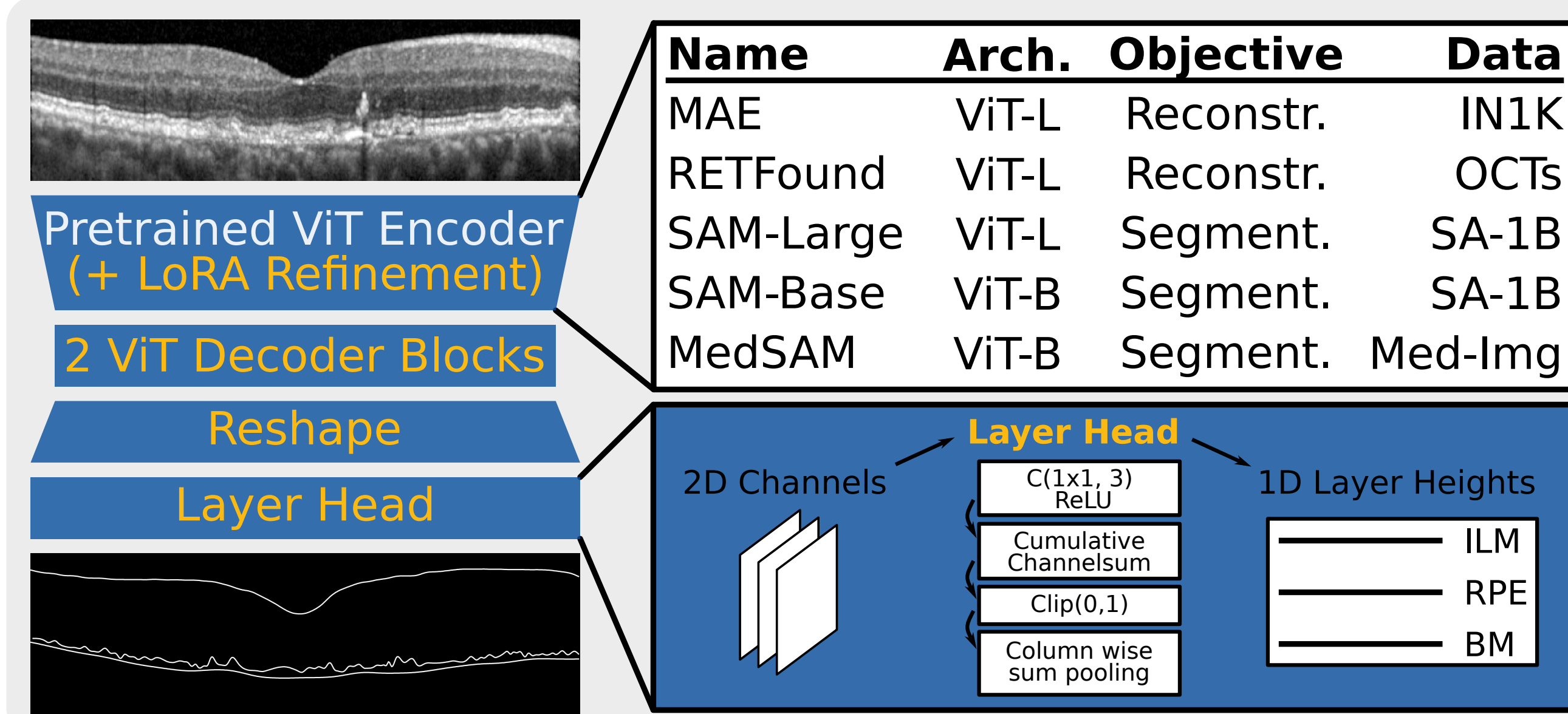
## Problem:

- Segmentation of retinal layers is critical in the analysis of optical coherence tomography (OCT) data.
- The appearance and quality of OCT datasets vary considerably.
- Our model for order-constrained regression<sup>1</sup> provides state-of-the-art results in the training domain but generalises poorly to unseen domains.

## Research Question:

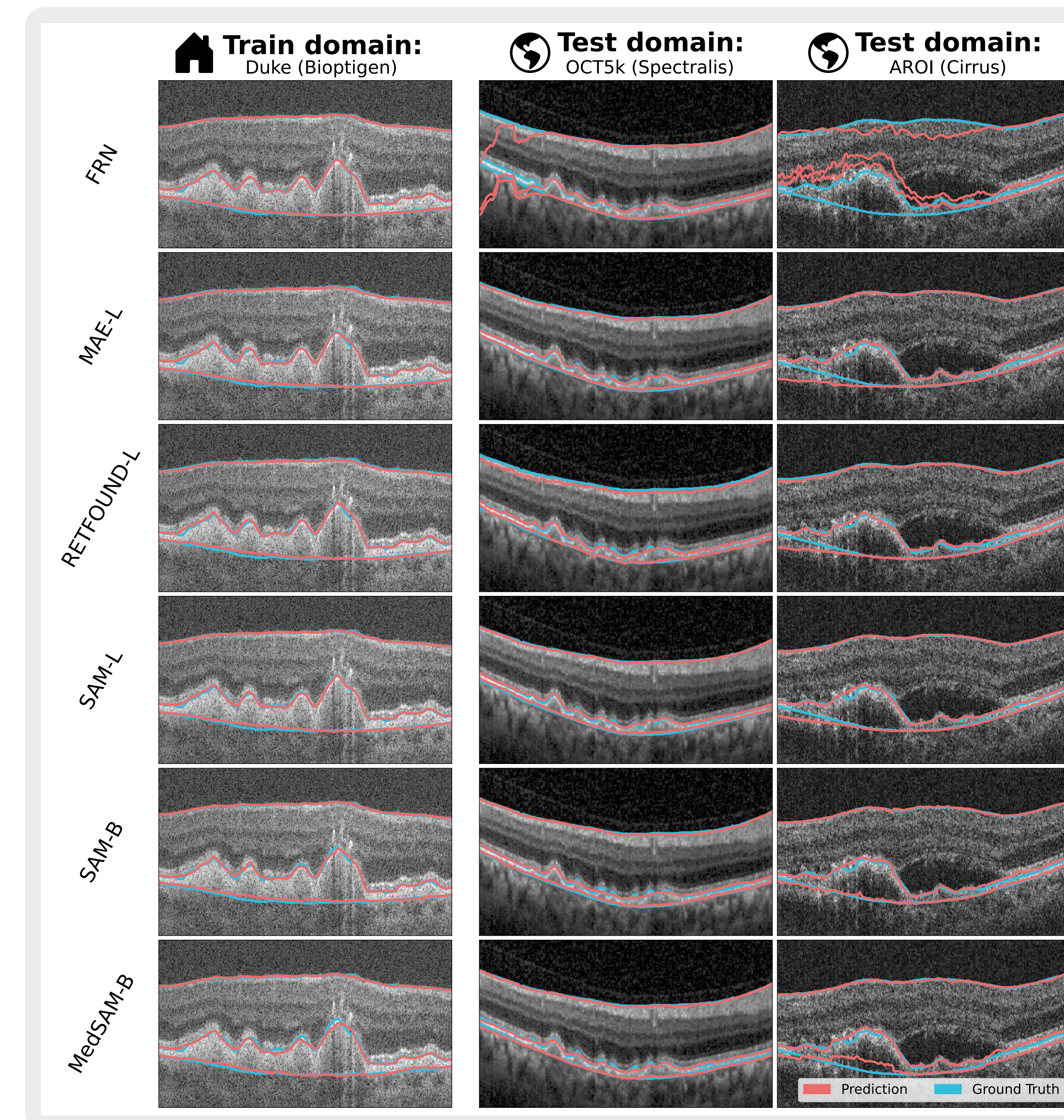
- Do vision foundation models improve the robustness of retinal layer segmentation models?
- Which foundation models work best?

## Methods:



- We compare ViT encoders of varying sizes, pre-trained with different data and objectives.
- Encoder parameters are fixed except for a LoRA<sup>9</sup> refinement
- 2 Transformer blocks and a layer head are trained to map ViT embeddings to layer heights.

## Results:



### In Domain:

- The CNN trained from scratch is the best by a small margin.

### Under Domain Shift:

- All tested pre-trained models do better on new domains than the baseline CNN. SAM-L is the best.
- RETFound-L domain specific pre-training reduces the performance compared to the more general MAE-L
- RETFound-L performance varies between test datasets. Is there a bias in the training data?

Comparison of the absolute segmentation error in pixel for the tested models on the training domain (Duke) and two separate datasets. We report the Median and 95% Quantile over all layers within each dataset.

Model	Duke <sup>2</sup> 🏠	OCT5k <sup>3</sup> 🌐	AROI <sup>4</sup> 🌐
CNN <sup>1</sup>	<b>0.21 (0.97)</b>	0.71 (44.62)	0.99 (21.75)
MAE-L <sup>5</sup>	0.23 (1.01)	0.37 (1.2)	0.37 ( <b>3.58</b> )
RETFound-L <sup>6</sup>	0.24 (1.02)	0.8 (2.73)	0.5 (4.68)
SAM-L <sup>7</sup>	0.22 (0.99)	<b>0.34 (1.1)</b>	<b>0.35 (3.88)</b>
SAM-B <sup>7</sup>	0.25 (1.18)	0.44 (2.63)	0.64 (10.37)
MedSAM-B <sup>8</sup>	0.25 (1.17)	0.43 (2.01)	0.7 (9.19)

## Conclusion

- ViT-based foundation models can increase the robustness of OCT layer segmentation compared to CNNs trained from scratch.
- The choice of the right foundation model matters.
- Within the foundation models, bigger models are better.
- Sometimes generic models outperform specialized models.
- LoRA refinement of the encoder helps all models, but does not change their relative performance.

### References

- Morelle et al. (2023) Scientific Reports 13 (1): 8162. <https://doi.org/10.1038/s41598-023-35230-4>.
- Farsiu et al. (2014) Ophthalmology 121 (1): 162–72. <https://doi.org/10.1016/j.ophtha.2013.07.013>.
- Arikan et al. (2023) bioRxiv. <https://doi.org/10.1101/2023.03.29.534704>.
- Melinščak, et al. (2021) Automatika 62 (3–4): 375–85. <https://doi.org/10.1080/00051144.2021.1973298>.
- He et al. (2022) CVPR 15979–88. <https://doi.org/10.1109/CVPR52688.2022.01553>.
- Zhou et al. (2023) Nature 622 (7981): 156–63. <https://doi.org/10.1038/s41586-023-06555-x>.
- Kirillov et al. (2023) ICCV 3992–4003. <https://doi.org/10.1109/ICCV51070.2023.00371>.
- Ma et al. (2024) Nature Communications 15 (1): 654. <https://doi.org/10.1038/s41467-024-44824-z>.
- Hu et al. (2022) ICLR <https://doi.org/10.48550/arXiv.2106.09685>