

Exploring finetuning ViT MAEs for fg/bg segmentation

[Klára Janoušková](#)

[Reconstruction finetuning on pascal:](#)

[Reconstruction finetuning on Berkely Deep Drive](#)

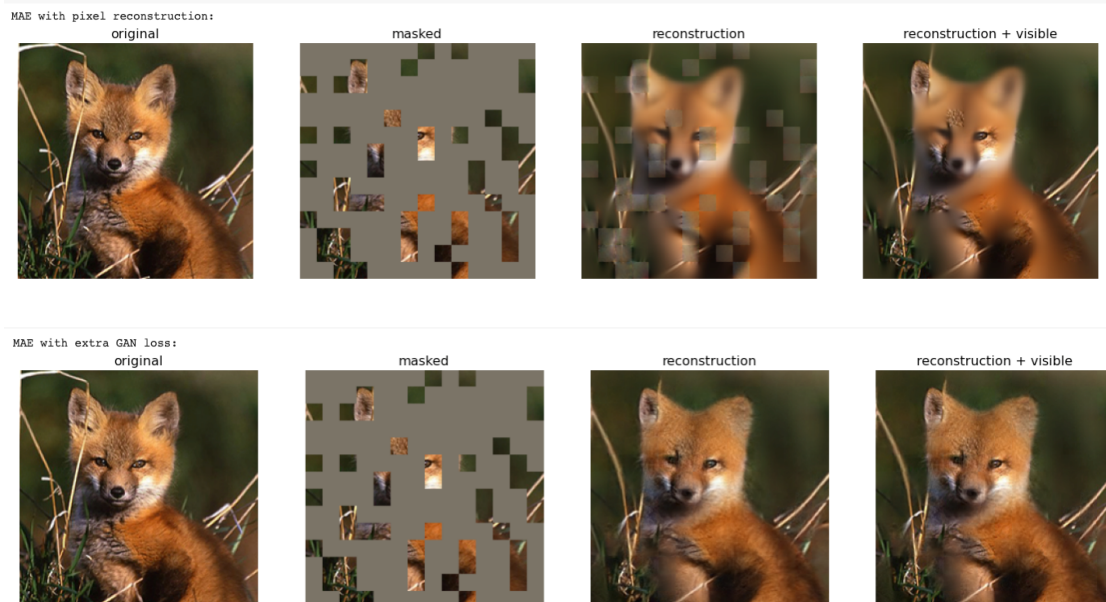
[Binary segmentation training](#)

[Reconstruction x segmentation loss interaction according to weight](#)

[Increased resolution to 384x384](#)

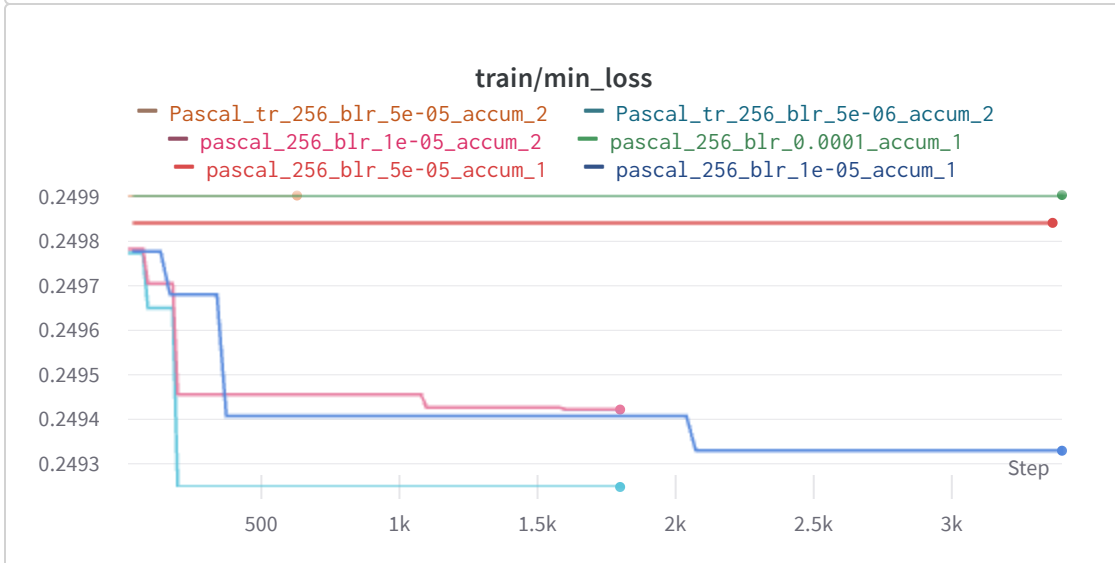
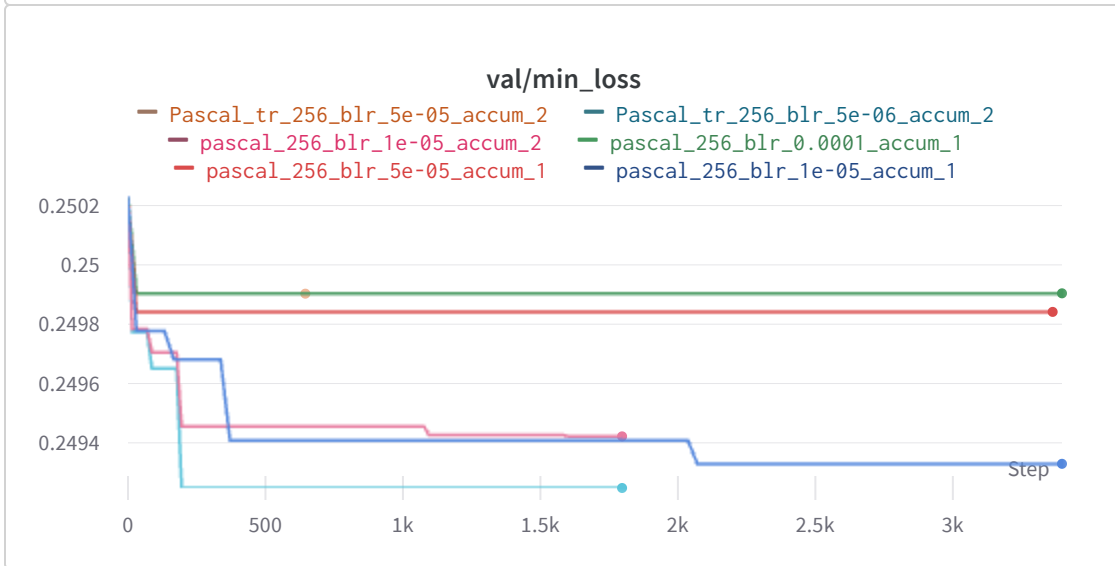
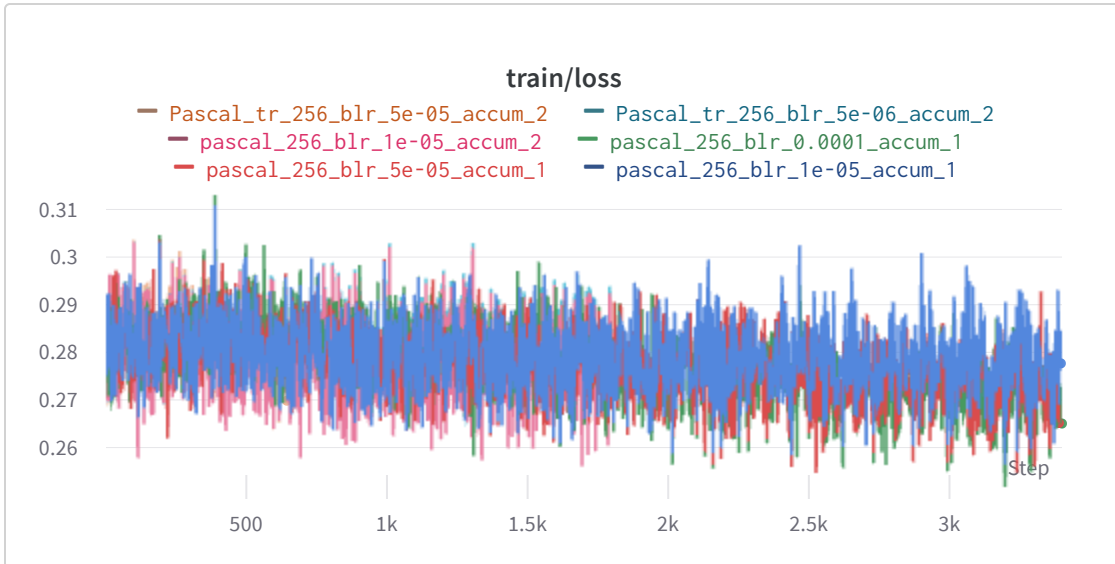
[Basic Sweep with ViT + ConvNet](#)

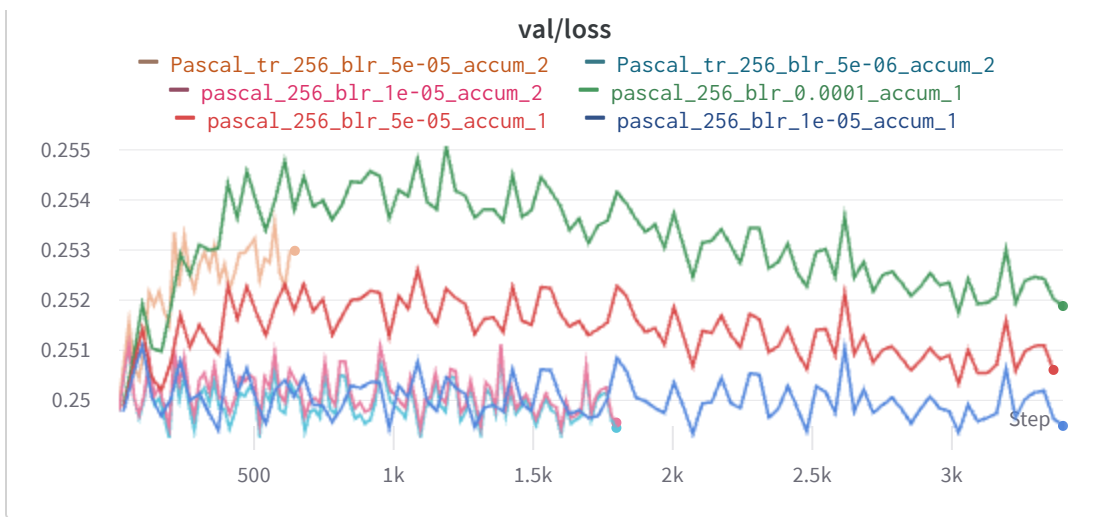
An example of the pretrained MAE model output from the authors, both with and without using GAN loss:



▼ Reconstruction finetuning on pascal:

The min val loss doesn't really change. If we increase lr, it starts increasing (not sure it is visible here). This could be because the pascal dataset is very similar to image-net.





▼ Reconstruction finetuning on Berkely Deep Drive

In these experiments, we try to change two things:

- 1) We chose a dataset that should be different enough from ImageNet, where the model was pre-trained
- 2) We remove optimization tricks from the finetuning code (layer decay, cosine scheduler, ...) and only use basic AdamW optimization with constant learning rate.

When we change the dataset, the training curves look good on the training dataset and they do decrease a bit on the validation set in the first few iterations. Still not ideal but better than before:





val/min_loss

Showing first 10 runs



train/loss

Showing first 10 runs



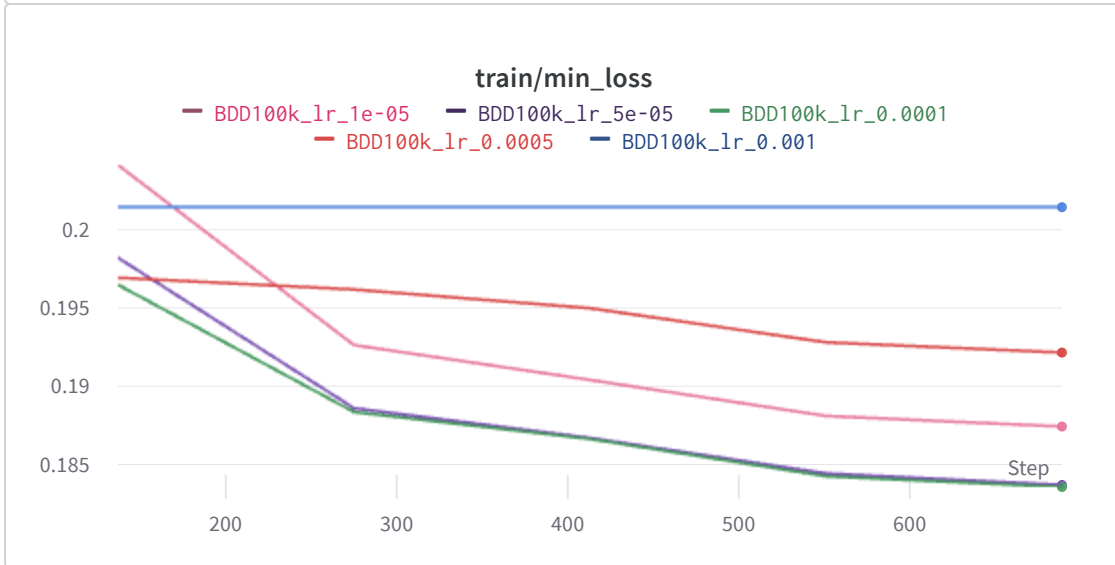
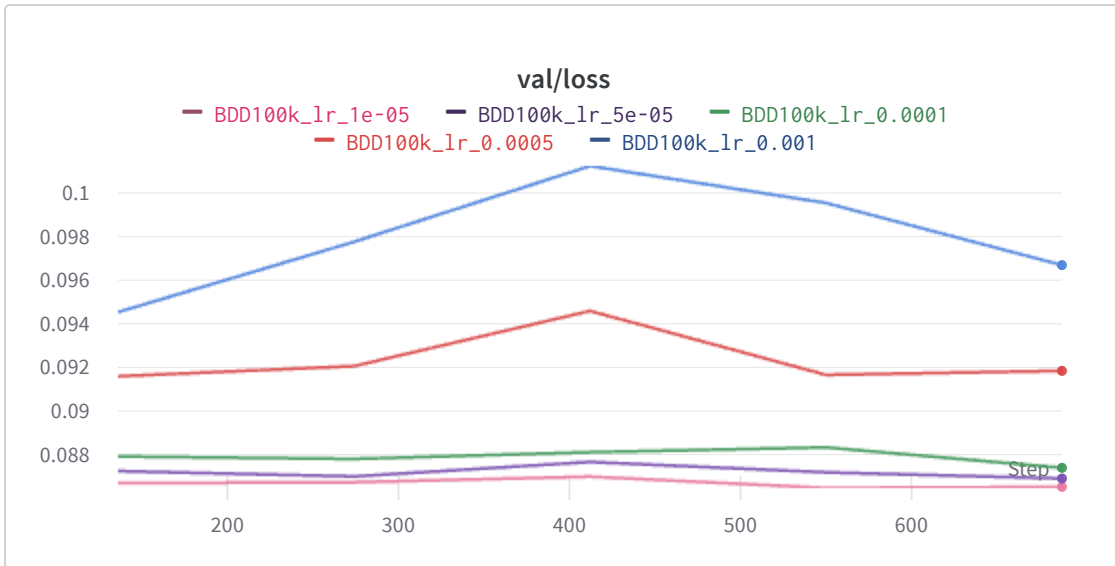
val/loss

Showing first 10 runs



After simplifying the optimization process, we get the following results:

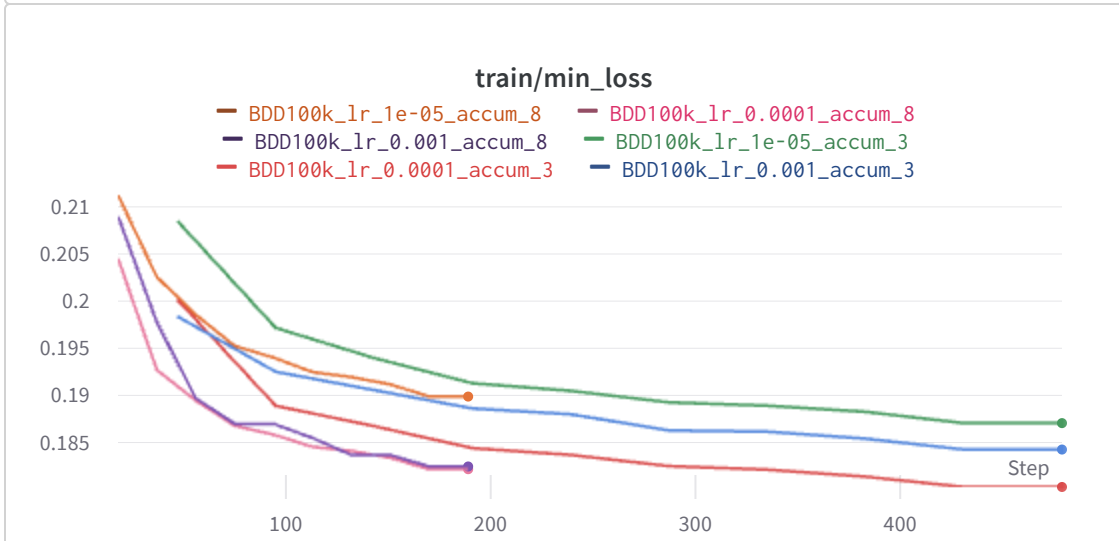
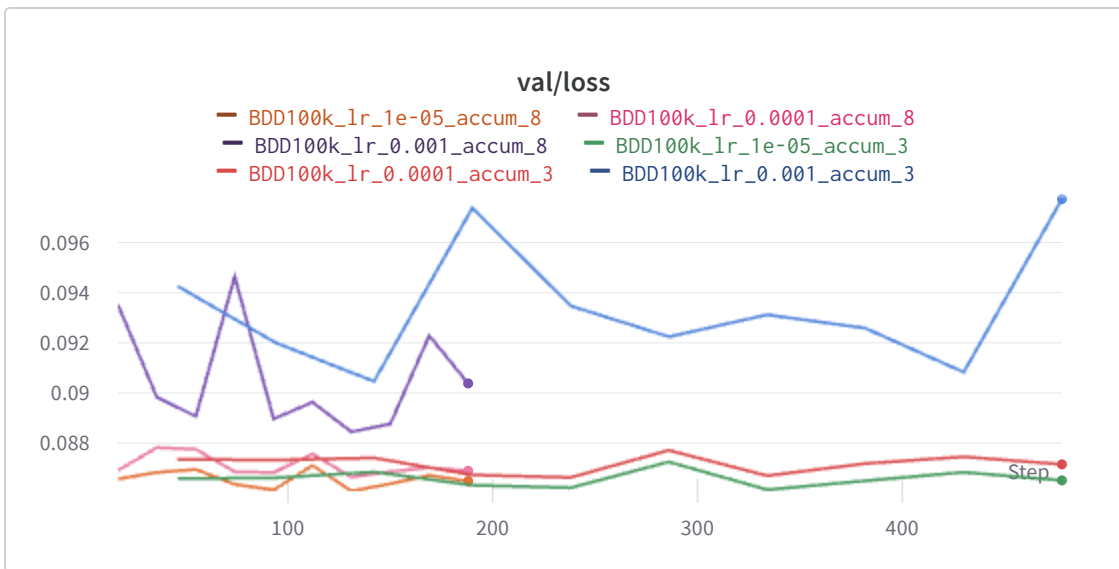
We see that as long as we keep the learning rate low enough, we see a healthy behaviour, if we increase it, the loss rather starts increasing on the validation dataset.

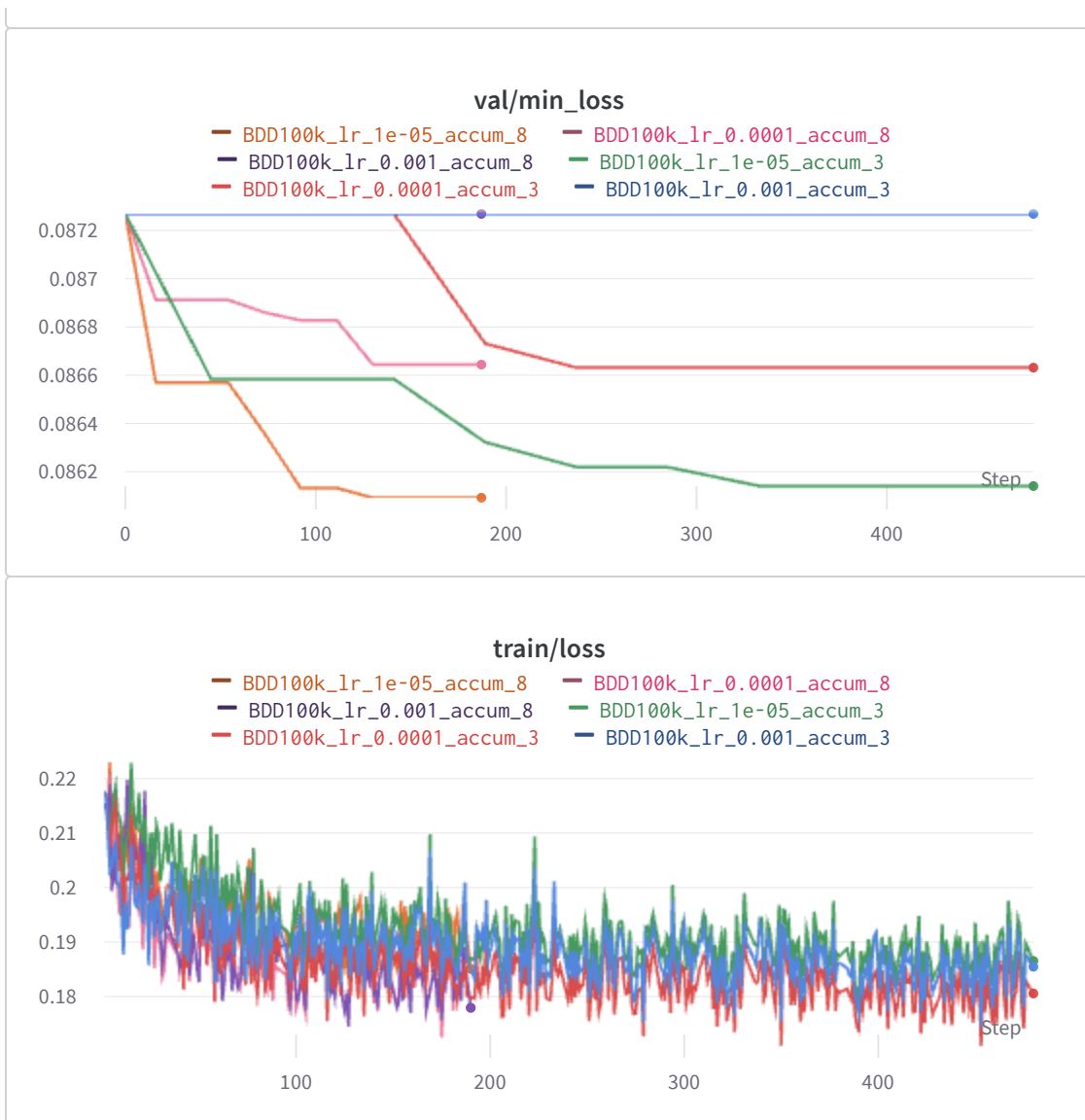




Playing around with more gradient accumulation steps (effective increasing batch size):

Same results as before, we see the smaller learning rate is important





TODO: Add images comparing pretrained/finetuned/gan models

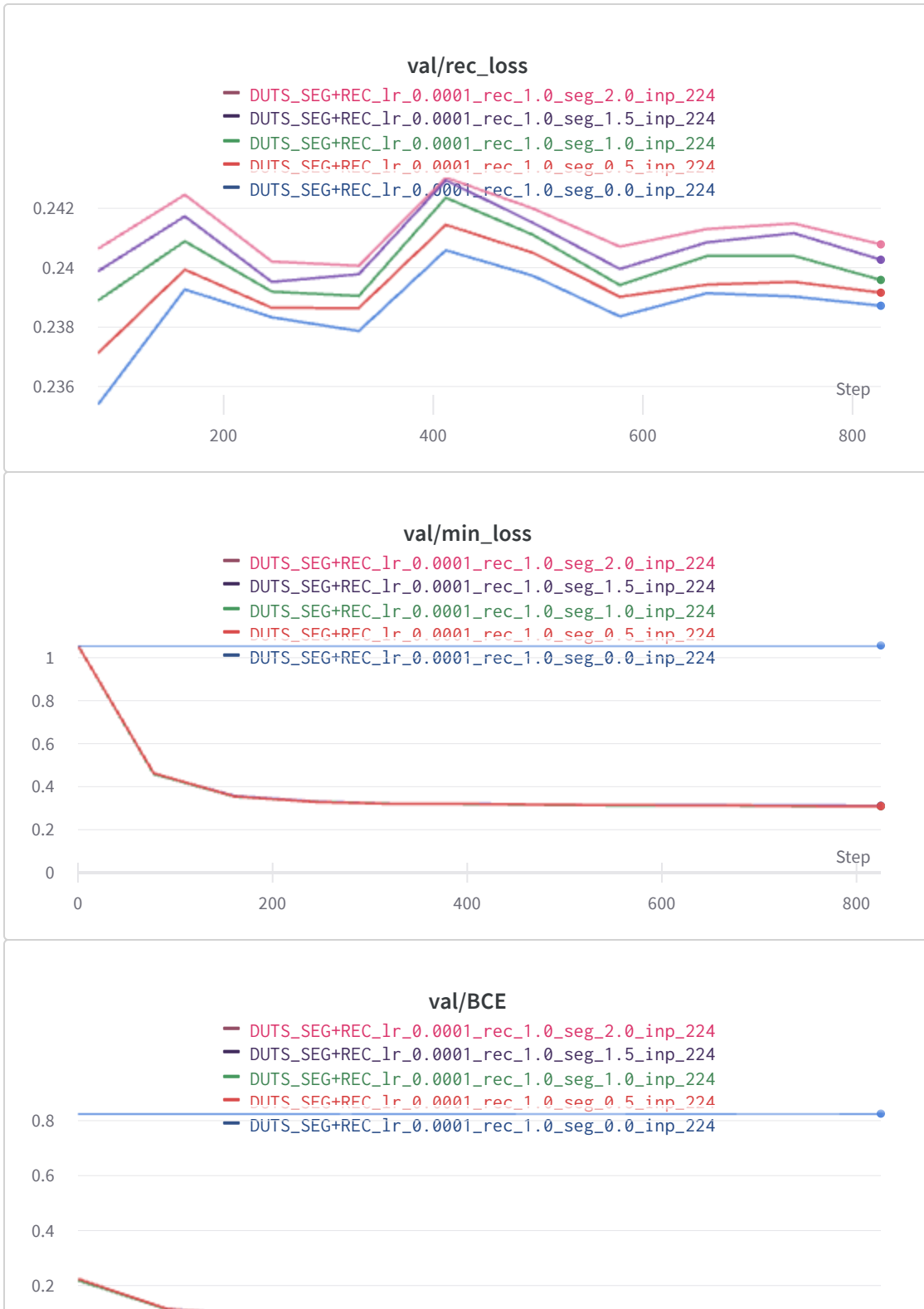
▼ Binary segmentation training

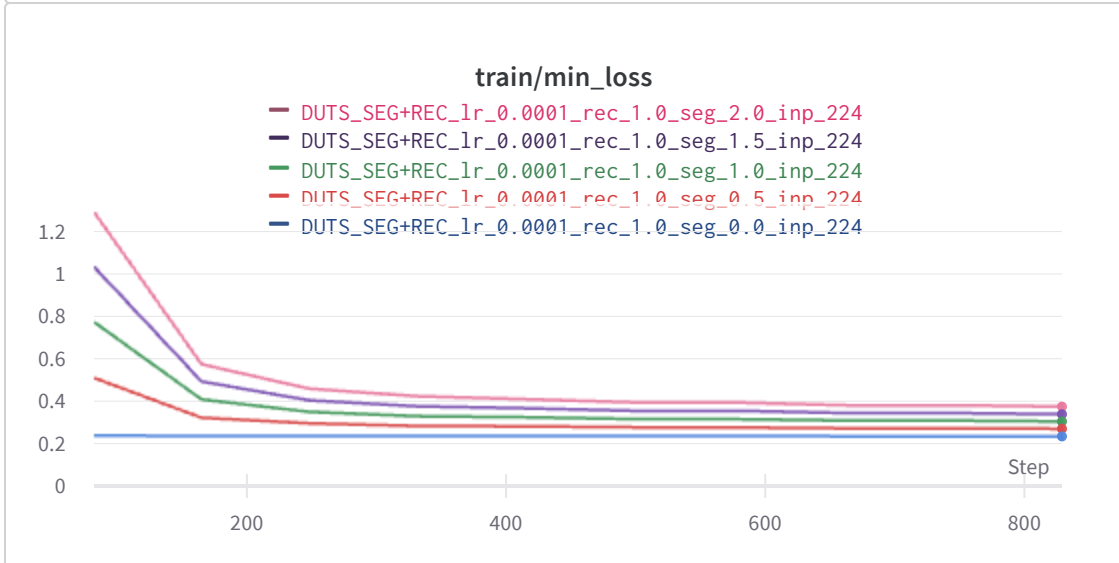
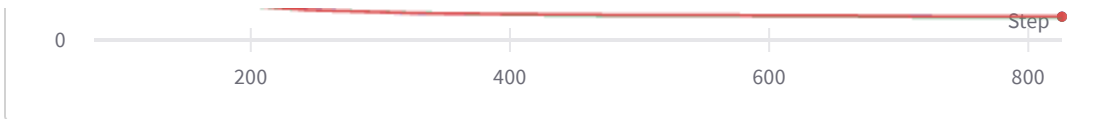
All the experiments here are on the DUTS dataset, which is based on image-net.

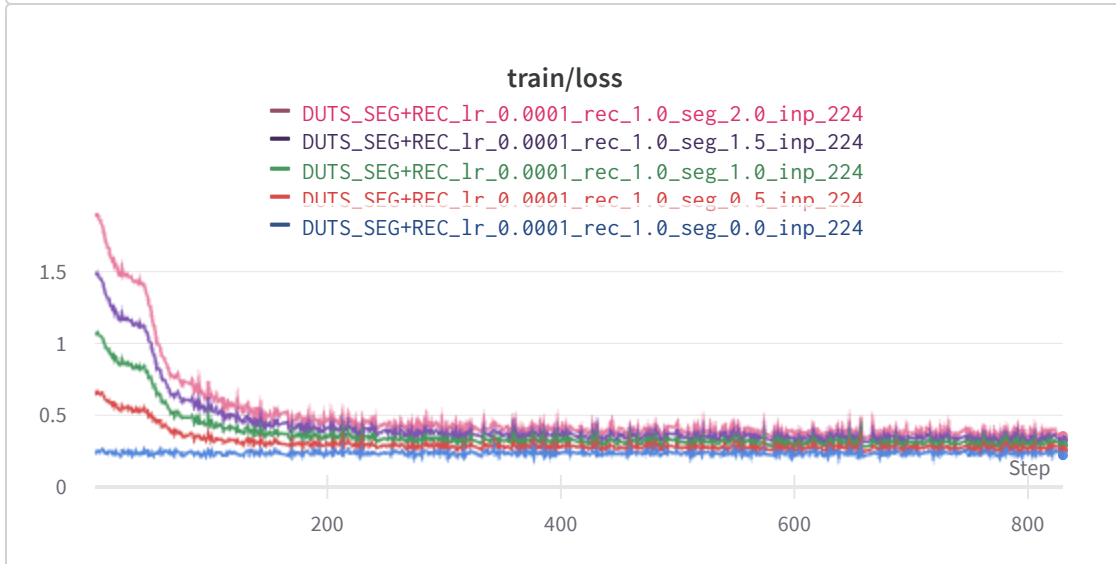
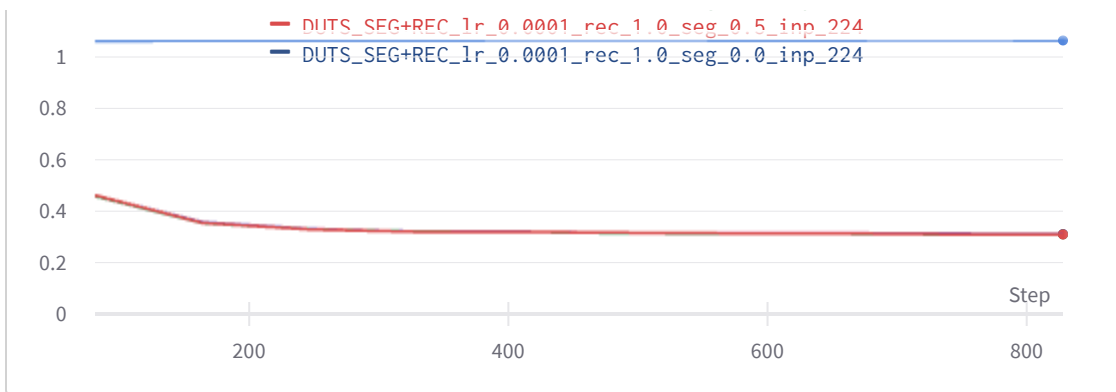
We are mostly exploring the qualitative results and how the two tasks interact.

▼ Reconstruction x segmentation loss interaction according to weight

We can see that the weight doesn't really have a big effect, the reconstruction is mostly constant while the segmentation loss is always decreasing at the same rate, as long as its weight is non-zero.







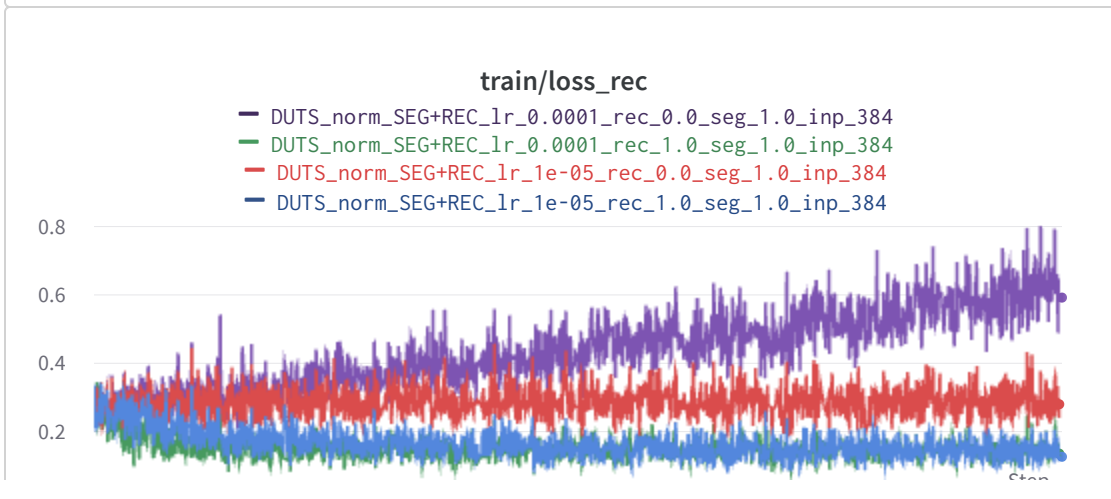
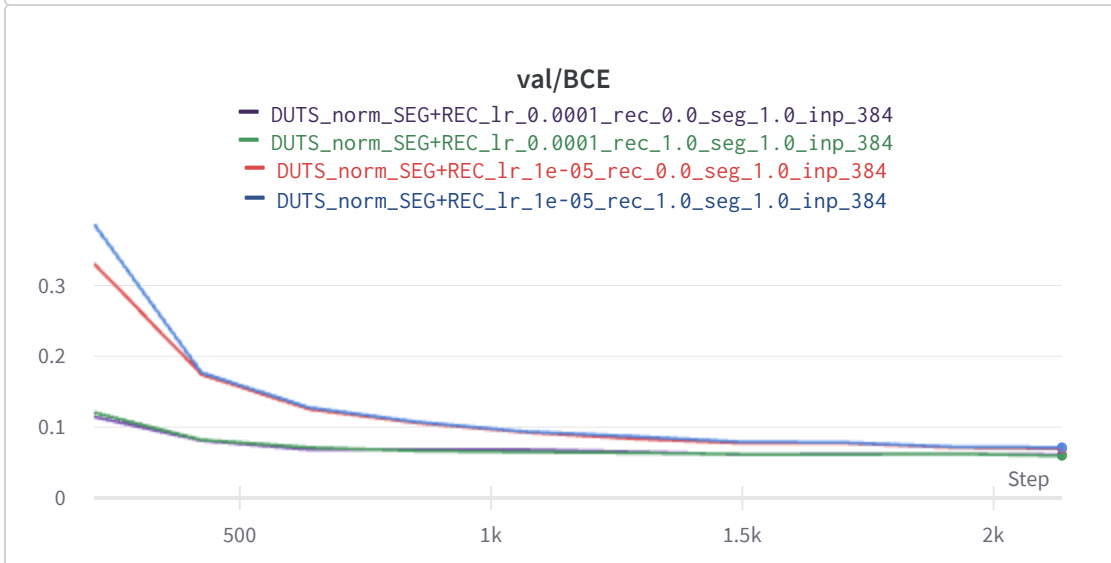
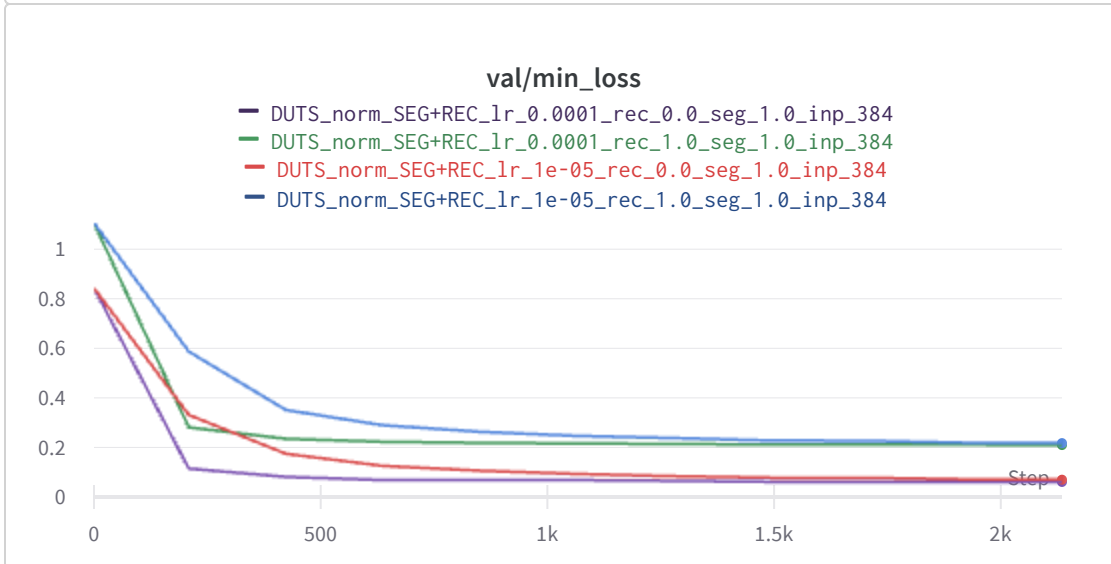
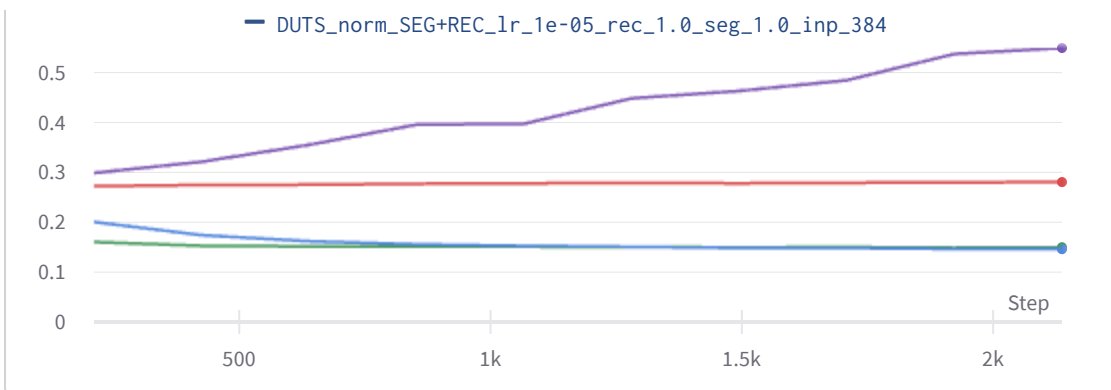
▼ Increased resolution to 384x384

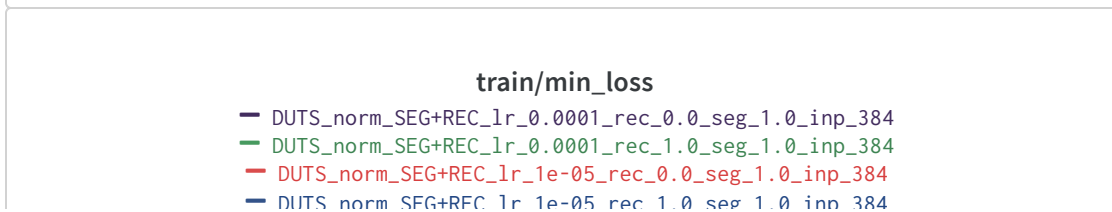
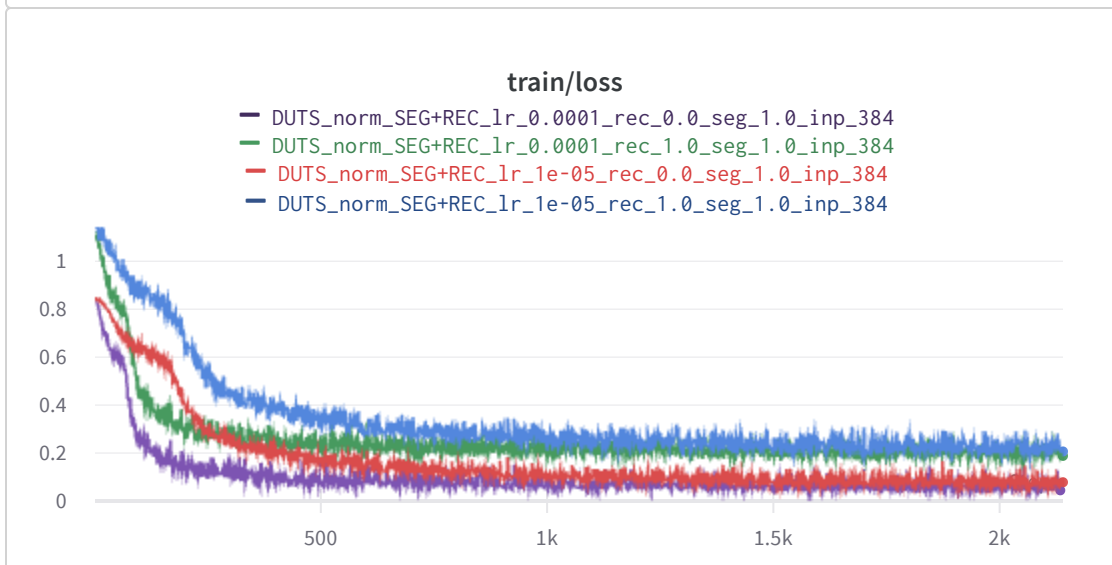
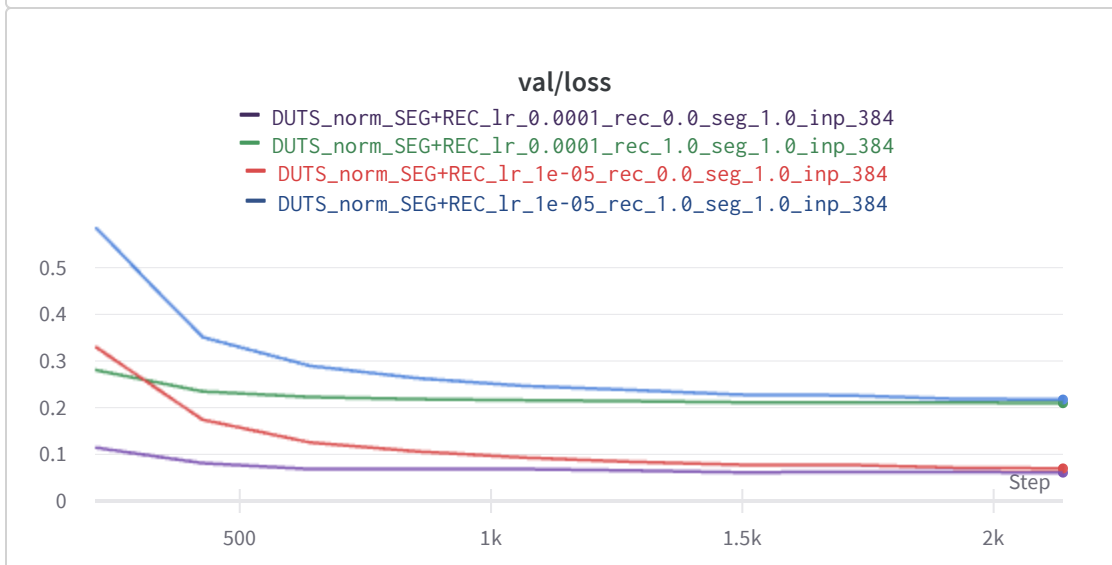
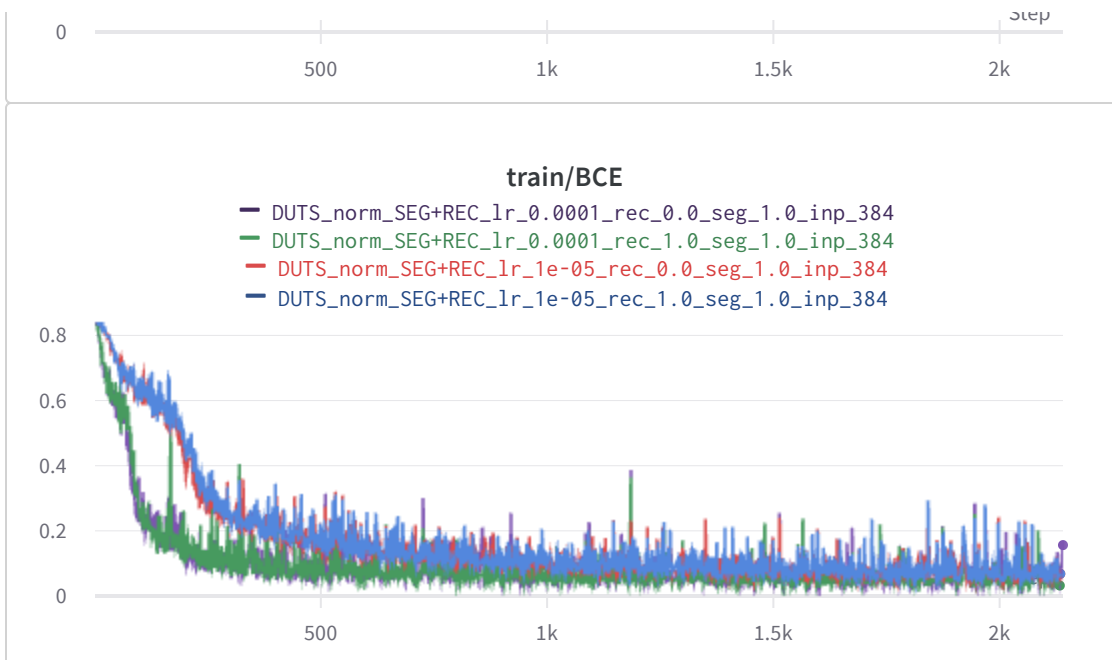
Results when increasing the resolution to 384x384:

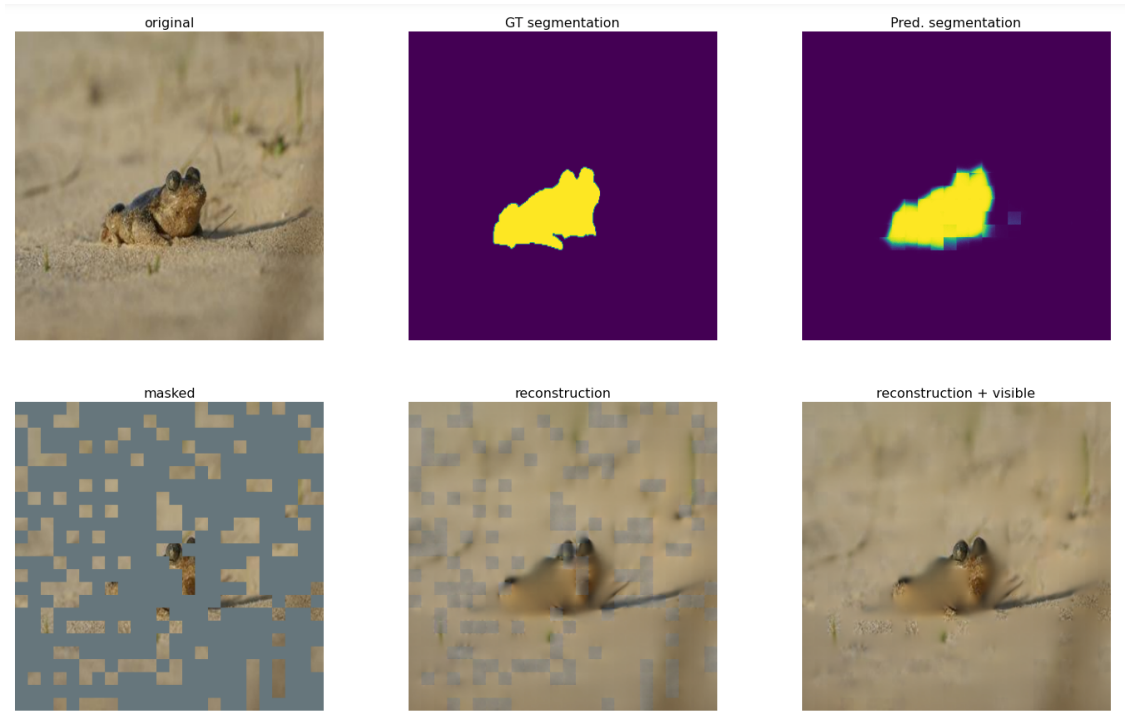
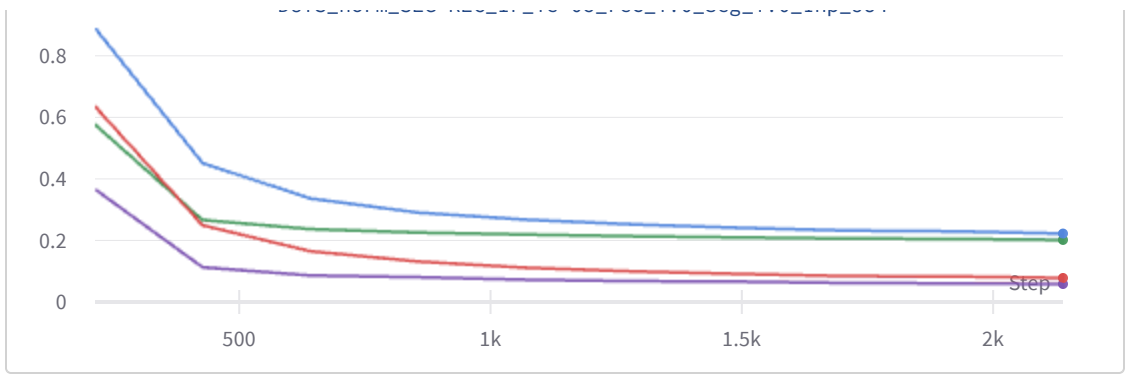
We see the results are a bit better thanks to the increased resolution, but we still have the issues with patch inconsistency, we will need to do something about the model architecture.

We can actually see the patch boundaries even in the reconstruction output, but it is not so visible as it is not binary.



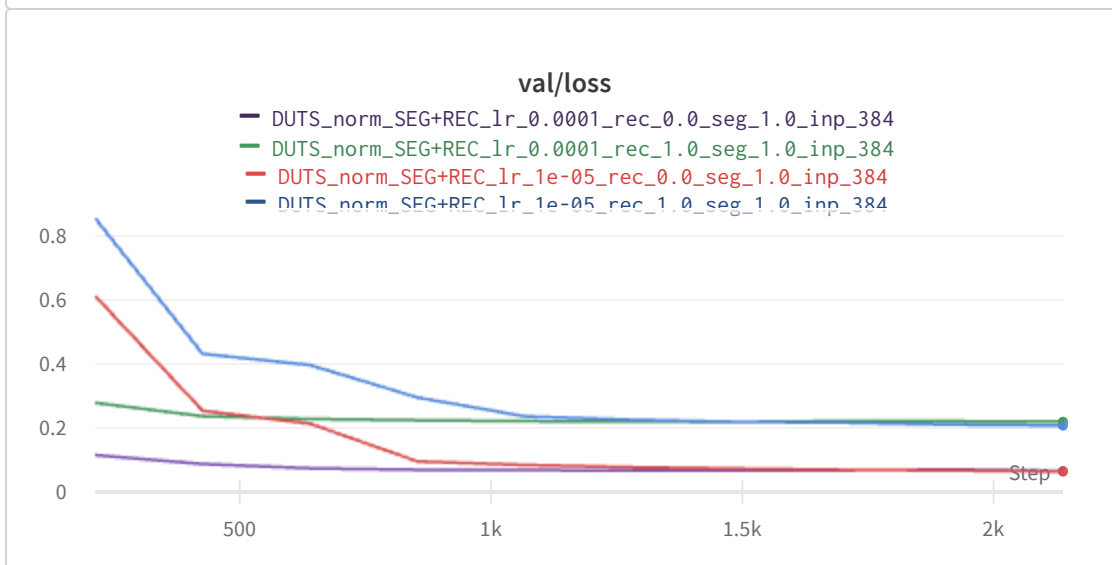
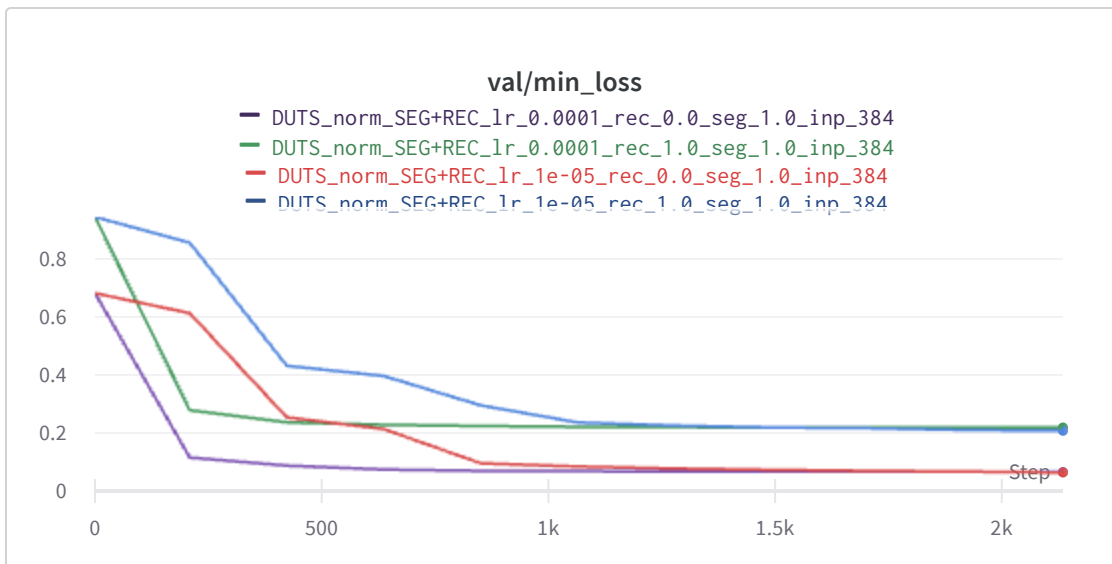


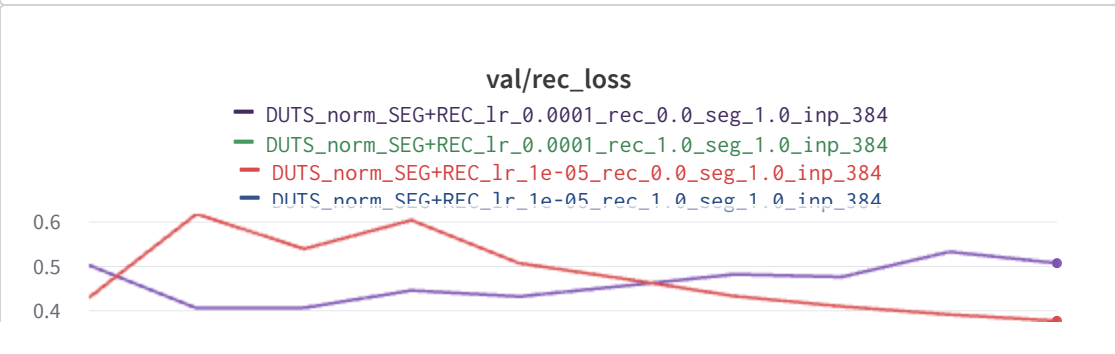
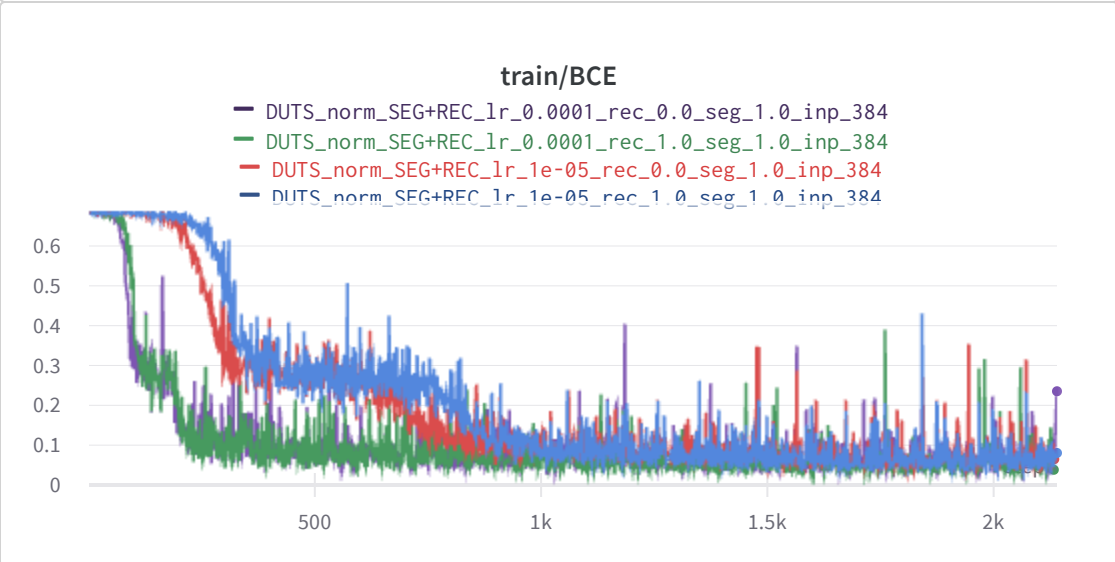
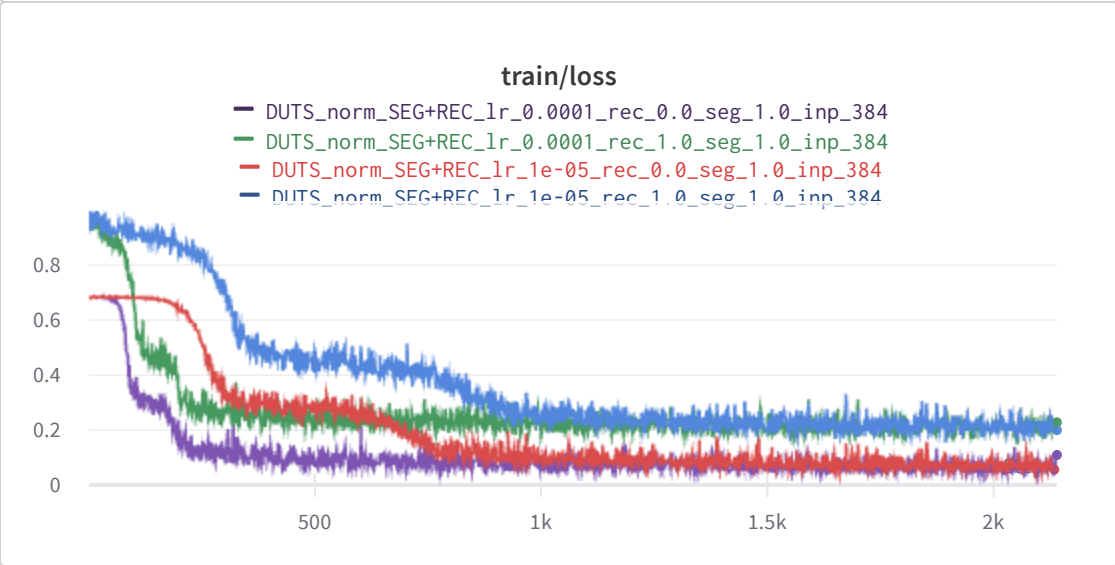
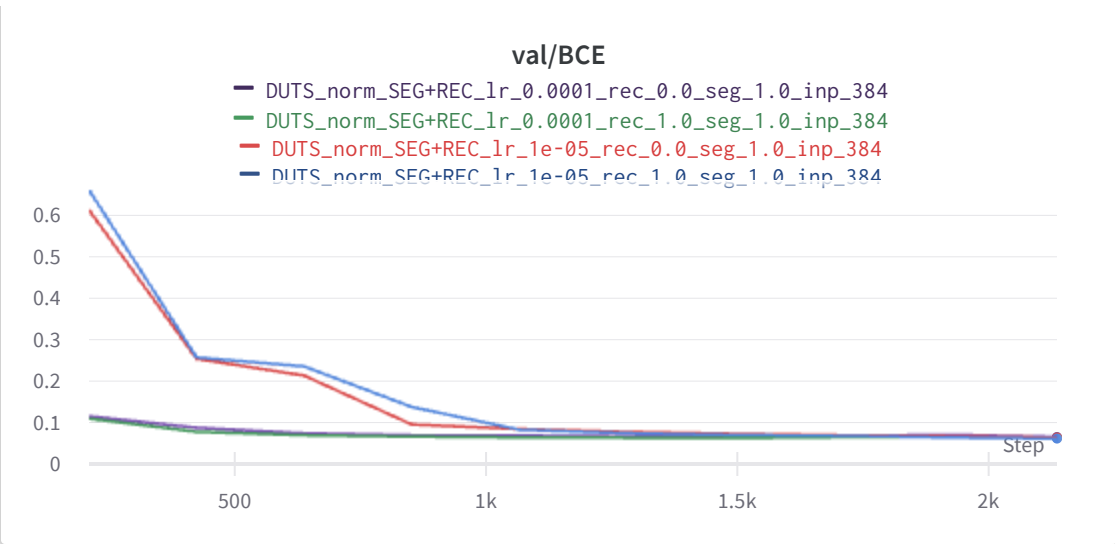


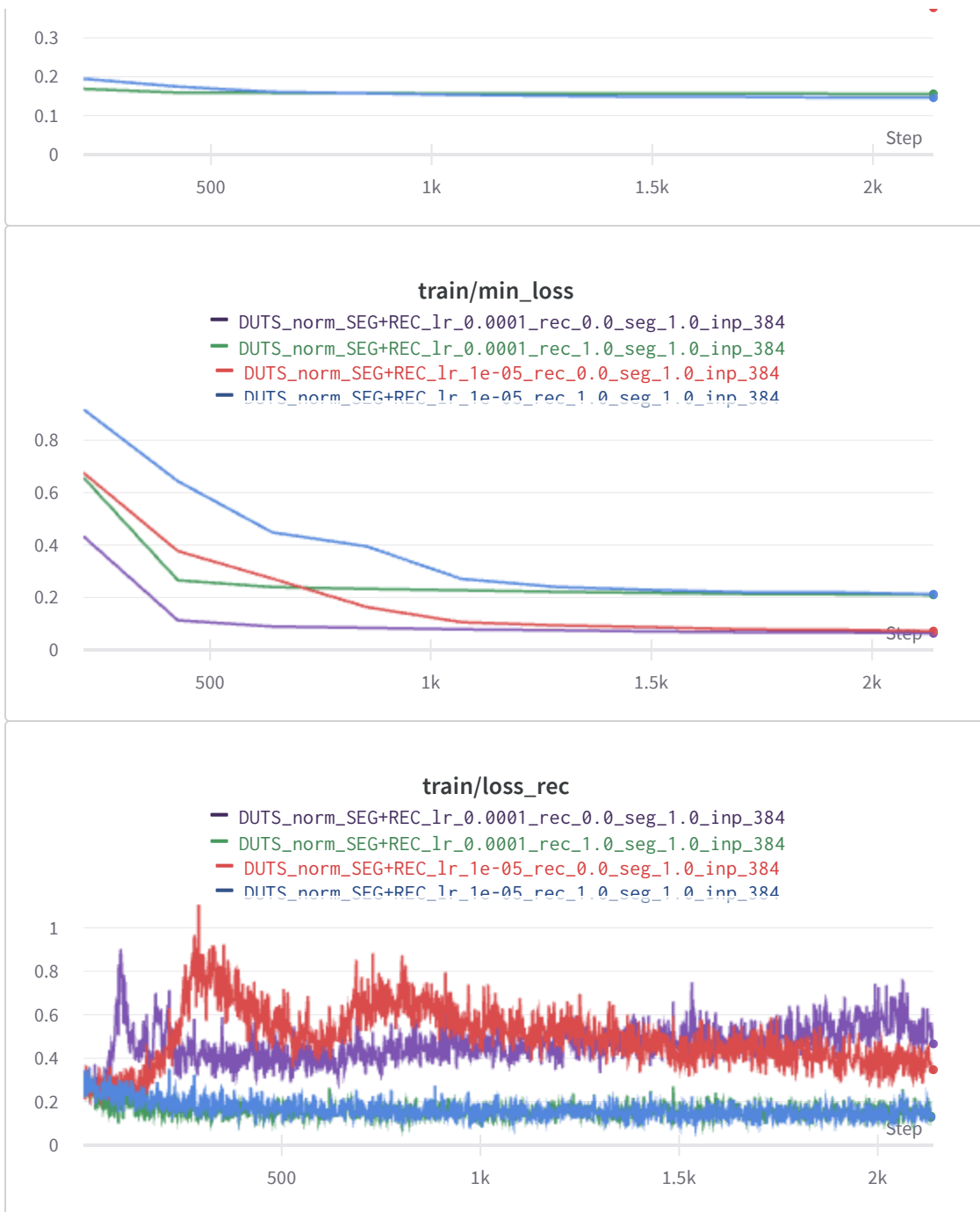




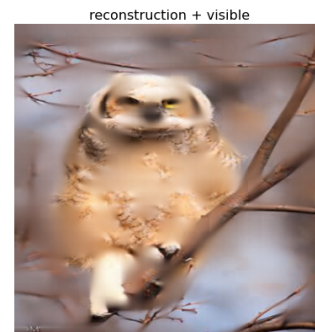
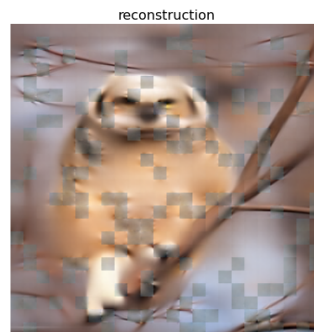
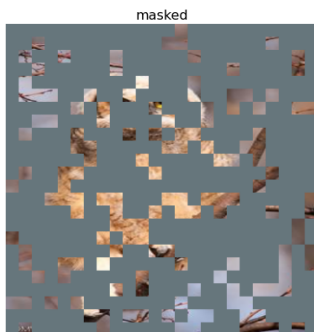
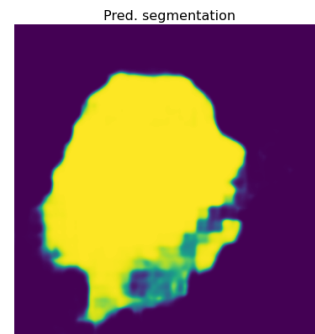
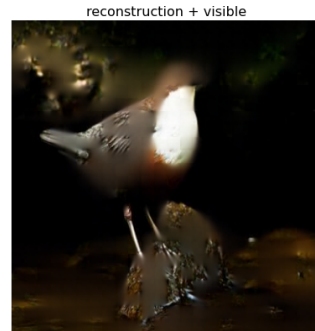
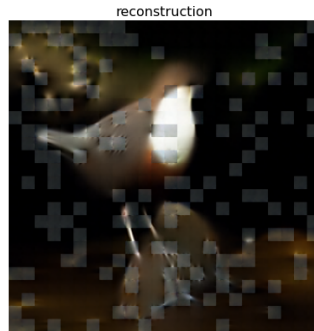
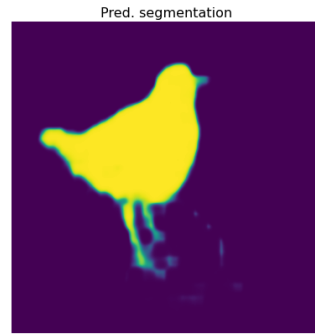
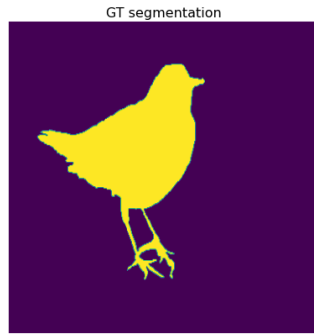
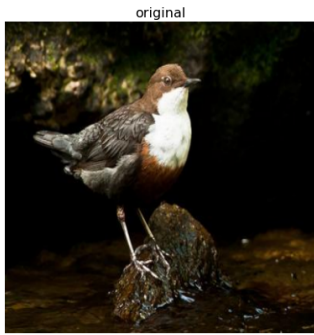
▼ Basic Sweep with ViT + ConvNet







We can still see some artifacts due to patches, but it has improved a lot. We might want to try a deeper convnet (should help with performance, the transformer decoder was much bigger), but most likely, we need to add a skip connection from the input image to the convnet part.



Created with  on Weights & Biases.

<https://wandb.ai/klara/MAE-finetune/reports/Exploring-finetuning-ViT-MAEs-for-fg-bg-segmentation--VmldzoyOTY5NjM4>

