

6. Supplementary of Semi-Supervised Semantic Segmentation of Vessel Images using Leaking Perturbations

6.1. Algorithm summary of the proposed model

Our model is summarized in Algo. 1. Furthermore, LeakGAN can be used for the cross-domain scenarios, which the unlabelled target \mathbf{X}_{ul}^t also is fed into the model, as shown in Algo. 2.

Algorithm 1: Training of the LeakGAN

Input: Source: $\mathbf{X}_l^s, y_l^s, \mathbf{X}_{ul}^s$
Result: Segmentation classifier (discriminator)

- 1 $\theta_D, \theta_G, \theta_{MT} \leftarrow$ initialization
- 2 **for** iterations of training **do**
- 3 $\mathbf{X}_l^s, y_l^s, \mathbf{X}_{ul}^s \leftarrow$ sample mini-batch
- 4 $z \leftarrow$ sample from $\mathcal{N}(0, I)$
- 5 Generate $\tilde{\mathbf{X}}_f$ by feeding z through generator G
- 6 $\mathcal{L}_{sup}, \mathcal{L}_{unSUP}, \mathcal{L}_{cons-FL}, \mathcal{L}_{adv-} \leftarrow$ calculated by Eq. (2), Eq. (3), Eq. (13) and Eq. (15)
- 7 **for** iterations of inference model updating **do**
- 8 $\theta_D \leftarrow -\Delta_{\theta_D} \mathcal{L}^*$
- 9 **end**
- 10 **for** iterations of discriminator updating **do**
- 11 $\theta_G \leftarrow -\Delta_{\theta_G} \mathcal{L}_{adv-}$
- 12 **end**
- 13 **end**

Algorithm 2: Training of the cross-domain LeakGAN

Input: Source: $\mathbf{X}_l^s, y_l^s, \mathbf{X}_{ul}^s$, Target: \mathbf{X}_{ul}^t
Result: Segmentation classifier (discriminator)

- 1 $\theta_D, \theta_G, \theta_{MT} \leftarrow$ initialization
- 2 **for** iterations of training **do**
- 3 $\mathbf{X}_l^s, y_l^s, \mathbf{X}_{ul}^s, \mathbf{X}_{ul}^t \leftarrow$ sample mini-batch
- 4 $z \leftarrow$ sample from $\mathcal{N}(0, I)$
- 5 Generate $\tilde{\mathbf{X}}_f$ by feeding z through generator G
- 6 $\mathcal{L}_{sup}, \mathcal{L}_{unSUP}, \mathcal{L}_{cons-FL}, \mathcal{L}_{adv-} \leftarrow$ calculated by Eq. (2), Eq. (3), Eq. (13) and Eq. (15)
- 7 **for** iterations of inference model updating **do**
- 8 $\theta_D \leftarrow -\Delta_{\theta_D} \mathcal{L}^*$
- 9 **end**
- 10 **for** iterations of discriminator updating **do**
- 11 $\theta_G \leftarrow -\Delta_{\theta_G} \mathcal{L}_{adv-}$
- 12 **end**
- 13 **end**

6.2. Implementation

Our experiments are based on implementations using TensorFlow [1]. Our model is trained by the patches randomly cropped from the training images. The patch size is 64×64 . The structure of the generator is similar to that of DCGAN [31], which performs well on image generation tasks. A 1-D random normal noise with length 100 is fed to the generator. It passes through a fully connected layer, and the activation is reshaped to 8×8 . There are overall six convolution layers following the fully-connected layer. LeakyReLU is used as the activation function except for the last layer, where $\tanh()$ is used instead. Transposed convolution layers with stride 2 are utilized for the upsampling. Batch normalization is added between the convolutional layers. The intermediate outputs of the generator, which have the same frame size as the discriminator's, are used as the leaked information. The discriminator is of U-Net style and has a structure similar to the model proposed in [28]. The kernel of the consecutive convolution layers is 3×3 , and max-pooling with step 2 is utilized to downsample the image patches. After four stacked blocks of convolution and downsampling layers, the patches are encoded to the size of $4 \times 4 \times 512$. The decoder of the discriminator has a structure symmetric to the encoder. The encoder's output, also the intermediate outputs of the generator, are concatenated to the decoder for the further upsampling. During the experiments, the output of the first convolutional layer of the generator (features with size 8×8) is utilized for the leaking module. It spends around 1.69s, 1.72s, and 4.70s to generate each image prediction in patch-level for DRIVE, STARE, and CHASE_DB1, respectively, when our model is run on a single NVIDIA GTX1080 GPU.

6.3. Qualitative results for STARE and CHASE_DB1

In this section, more visualization of the experiments on STARE and CHASE_DB1 are demonstrated.

The whole-image prediction of the STARE and CHASE_DB1 are shown in Fig. 6.

Also, in the sixth column of Fig. 7a, the vessels of STARE are detected well though it has some blue pixels around them. However, some vessels' structures disappear (as false negative) in the predictions, e.g. as shown in the fourth column of Fig. 7b.

6.4. More scenarios of cross-domain evaluations

We extend the cross-domain segmentation evaluation onto another scenario: training the LeakGAN framework on the healthy images and testing it on the pathological images. For the DRIVE dataset, there are seven images from diabetic patients. We also evaluate the performance for the health retina images as training images, the diabetic patient's retina images for the test. The results are given in Table 7, and

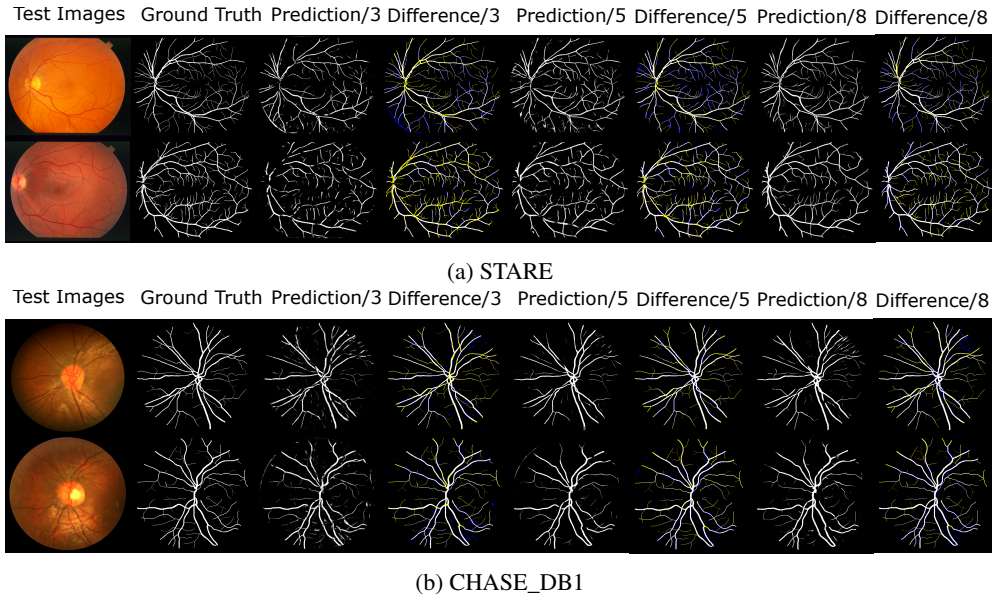


Figure 6: Visualization of the semantic segmentation on testing images. The first two columns are test images and their ground truth. Then the results are shown when 3, 5 and 8 labelled training images are used. Its differences to ground truth are shown next to the prediction columns, with blue pixels being false positive, and yellow being false negative.

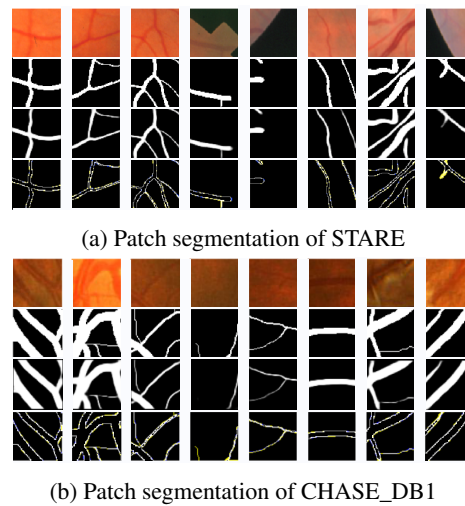


Figure 7: Patched segmentation results for the scenario of 8 labelled training image using different datasets. The testing patches and their ground truth are in first and second row respectively. The third row is for the predictions. The difference between the prediction and ground truth is shown in the last row, blue pixels being false positive, and yellow being false negative.

segmentation examples are shown in Figure 8a. Also the visualization of the scenario from STARE to DRIVE is shown in Fig. 8b. From these results, we can see that the model performs reasonably in this scenario.

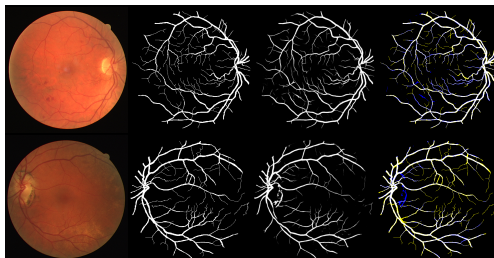
Furthermore, we evaluate our model’s generalization when it is applied to unseen datasets. During the experiments, the model is trained following Algorithm 1 on one dataset,

while the testing is done on retina images of another dataset. For example, these results are obtained: $Acc(\%)=95.06$, $Sp(\%)=92.15$, and $Se(\%)=83.37$, when the transfer from DRIVE (training) to STARE (testing) is considered. This shows that our model also generalizes well to unseen data resources.

Table 7: Performance for training on DRIVE(H), testing on DRIVE(D). Numbers for the training setup are for labelled and unlabelled images respectively.

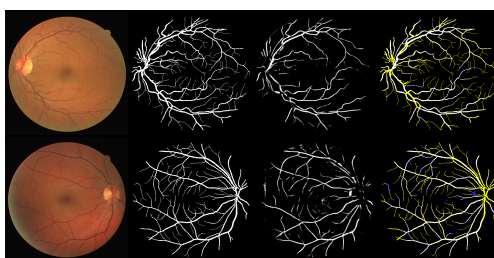
Train Dataset Test Dataset	DRIVE(H)(5+5) DRIVE(D)	DRIVE(H)(10+10) DRIVE(D)	DRIVE(H)(10+23) DRIVE(D)
Acc(%)	92.15	92.23	93.61
Sp(%)	85.39	85.32	80.52
Se(%)	97.53	95.16	98.30

Test Images Ground Truth Prediction/8 Difference/8



(a) From DRIVE (Healthy) to DRIVE (Diabetic)

Test Images Ground Truth Prediction/8 Difference/8



(b) From STARE to DRIVE

Figure 8: Visualization of the cross-domain prediction results. The test images and their ground truth are shown in the first two columns. The predictions are listed on the third. The difference between the prediction and ground-truth is shown in the last column overlaying on ground-truth, with blue pixels indicating false positive, and yellow indicating false negative.