

## Supplementary Material

In this supplementary material we give a detailed overview of the training and evaluation settings along with hyperparameters. We also provide a detailed description of the employed augmentation methods. Further, we show additional ablation studies and investigations. We also depict the SegFormer architecture in a block diagram. Finally, we discuss limitations and ethical implications for our method.

### A. Detailed Description of Augmentation Methods

In the following section, we present a detailed description of the employed augmentation methods.

**Random crop:** Parameter  $\rho$  defines the maximum proportion a single class can occupy in the random crop.

**PhotoAug:** Photometric augmentation (PhotoAug) comprises the following steps: Each transformation is applied to the image with a probability of 0.5. The position of the random contrast adjustment is in second (mode ①) or second to last position (mode ②). The position is randomly selected for each image.

1. random brightness
2. if ①: random contrast
3. convert color from RGB to HSV
4. random saturation
5. random hue
6. convert color from HSV to RGB
7. if ②: random contrast
8. randomly swap channels

**Bilateral filter:** A bilateral filter smoothes an image while preserving sharp edges. Its focus is on the removal of noise and textures. In general, it is a Gaussian filter that smoothes less in non-uniform regions (edge regions) and more in uniform image regions (non-edge regions). The bilateral filter at pixel index  $i$  can be described as:

$$G(i) = \frac{1}{w} \sum_{j \in \mathcal{J}} N(\Delta_{ij}; \sigma_s^2) N(\|\mathbf{x}_i - \mathbf{x}_j\|_2; \sigma_c^2) \mathbf{x}_j, \quad (3)$$

with

$$N(d; \sigma^2) = e^{-\frac{1}{2}(\frac{d}{\sigma})^2}, \quad (4)$$

and normalizing factor

$$w = \sum_{j \in \mathcal{J}} N(\Delta_{ij}; \sigma_s^2) N(\|\mathbf{x}_i - \mathbf{x}_j\|_2; \sigma_c^2). \quad (5)$$

The neighboring pixel index is denoted as  $j$  and stems from the neighborhood  $\mathcal{J}$ . The neighborhood is defined by the kernel size which we sample from a uniform distribution between 1 (px) and 15 (px). Distances  $\Delta_{ij} = \sqrt{|h_i - h_j|^2 + |w_i - w_j|^2}$  and  $\|\mathbf{x}_i - \mathbf{x}_j\|_2$  denote the (Euclidean) pixel distance and color difference, respectively,

Listing 1. Pseudo-code of the PixMix [13] data augmentation.

```
def PixMix(x, z, K, beta): # mixing image z in Z
    x0 = random.choice({augment(x), x})
    # random number of mixing rounds
    for k = 1: random.choice({0, 1, ..., K}):
        xmix = random.choice({augment(x), z})
        mix_op = random.choice({add, multiply})
        xk = mix_op(xk-1, xmix, beta)
    return xk
```

where  $h_i$  and  $w_i$  are the height and width position of pixel  $i$  and  $\mathbf{x}_i \in \mathbb{G}^3$  denotes the vector of RGB values for pixel  $i$  (likewise for pixel  $j$ ). We set the spatial distance to  $\sigma_s = 75$  and the color distance to  $\sigma_c = 75$ . The probability for applying this filter is set to  $p = 0.5$ .

**PixMix:** The PixMix [13] augmentation method comprises multiple processing steps as shown in Listing 1. Here,  $\mathbf{x}$  denotes the input image and  $\mathbf{z} \in \mathcal{Z}$  denotes the mixing image from the PixMix set of fractal images  $\mathcal{Z}$  [13]. We set the maximum number of mixing rounds to  $K = 3$ . Note that the for loop is not executed for random choice = 0. The **mix\_op**( $\cdot$ ) function is randomly chosen to be either addition (**add**) or multiplication (**multiply**). It gets the images  $\tilde{\mathbf{x}}_{k-1}$  and  $\tilde{\mathbf{x}}_{\text{mix}}$  as inputs, as well as  $\beta = 3$ , which is used to generate independent weighting factors for the images. The weighting factors are sampled from a Beta distribution. For the **augment**( $\cdot$ ) function in the PixMix pseudocode we used only baseline augmentation methods PhotoAug and Random Flip (cf. Figure 3), that is why we denote the method as PixMix\*.

### B. Training/Evaluation Settings, Hyperparameters

In the following section, we will provide a detailed overview of the training and evaluation settings and hyperparameters. For the training and evaluation we employ PyTorch v.3.8.13 and the MMSegmentation toolbox v.0.11.0. Additionally, we refer to our repository, where all code for the conducted experiments is made available<sup>3</sup>.

**Training phase:** In Table 9 we list all settings and hyperparameters that were used for the training process. The polynomial learning rate schedule is defined as follows:

$$\eta(\tau) = \eta_0 \left(1 - \frac{\tau}{\tau_{\text{max}}}\right)^{0.9}, \quad (6)$$

with  $\eta(\tau)$  being the learning rate at optimizer step (iteration)  $\tau$  and  $\eta_0$  being the initial learning rate. The maximum number of iterations is given by  $\tau_{\text{max}}$ .

During training, the images from the source domain  $\mathcal{D}^S$  get resized to a resolution of  $720 \times 1280$ .

**Evaluation phase:** For evaluation, we always employ the final model weights after the full training and do not perform any checkpoint selection. We resize the input images

<sup>3</sup>Code is available at <https://github.com/ifnspaml/ReVT>

Table 9. **Settings and hyperparameters** for the SegFormer and DeepLabv3+ training.

Setting / Hyperparameter	SegFormer	DeepLabv3+
Optimizer	AdamW [24]	SGD
# of training iterations ( $\tau_{\max}$ )	40,000	60,000
Momentum values (AdamW) ( $\beta_1, \beta_2$ )	0.9, 0.999	-
Momentum ( $\beta$ )	-	0.9
Warm-up iterations	1500	-
Warm-up ratio	$1 \cdot 10^{-6}$	-
Initial LR ( $\eta_0$ )	$6 \cdot 10^{-5}$	$1 \cdot 10^{-3}$
Weight decay $\omega$	0.01	0.0005
Learning rate (LR) schedule ( $\eta(\tau)$ )	Polynomial (6)	Polynomial (6)
Batch size	2	2
Random decoder init	Kaiming initialization	Kaiming initialization
Resized input resolution $\mathcal{D}_{\text{train}}^{\text{GTA5}}$	$720 \times 1280$	$720 \times 1280$

Table 8. **Image and label resolution** [ $\text{px} \times \text{px}$ ] for the employed **evaluation** datasets. <sup>◦</sup>In both GTA and KITTI, there are images that differ by a few pixels from their normal resolution. <sup>\*</sup>The crowd-sourced Mapillary Vistas dataset does not have a fixed resolution, but a highly variable one.

Dataset name	Resolution ( $H \times W$ ) of ...	
	resized images	labels
GTA5 [32]	$512 \times 932^\circ$	$1052 \times 1914^\circ$
SYNTHIA [33] (SYN)	$512 \times 862$	$760 \times 1280$
Cityscapes [7] (CS)	$512 \times 1024$	$1024 \times 2048$
Mapillary Vistas [26] (MV)	various*	various*
BDD100k [43] (BDD)	$512 \times 910$	$720 \times 1280$
ACDC [34]	$512 \times 910$	$1080 \times 1920$
KITTI [1] (KIT)	$309 \times 1024^\circ$	$375 \times 1242^\circ$

during evaluation in a way that the image will be rescaled as large as possible within a pre-defined scale ( $512 \times 1024$ ), while still keeping their aspect ratios. The network output is then resized to the original image resolution. The mIoU is calculated on the original resolution, also referred to as label resolution. The resized image and the original label resolutions for all employed datasets are listed in Table 8.

The frame rate computations were performed on the rescaled Cityscapes dataset. We used 200 images for inference and computed the mean frame rate after a warmup phase of five images to account for any delays due to image reading operations.

### C. Additional Details on the Choice of Optimizer

In our experiments in Table 4 we show that the choice of optimizer has a strong effect on the baseline perfor-

mance of the models, as well as on the performance after re-parameterization. In Figure 5, we compare the SegFormer architecture with its standard optimizer setup (left) and the DeepLabv3+ architecture with its standard optimizer setup (right). We show the mean cosine similarity between three baseline (Ⓐ) models for the encoder only ( $\theta^E$ , upper plots) and the full model ( $\theta$ , center plots), during the training process. Note that the standard number of iterations differs for both models and is 40,000 for the SegFormer and 60,000 for the DeepLabv3+.

It can be seen that the mean cosine similarity for the network parts that are re-parameterized in our method (encoder only) have a similar mean cosine similarity after the training for both networks (0.995). In the bottom plots of the figure, we report the mIoU values for both in-domain (GTA5, green) and out-of-domain (OOD) data (Cityscapes, red) for the baseline (dashed lines) and re-parameterized (solid lines) models. It can be seen that the performance of the re-parameterized models is higher for the SegFormer for any training iteration. For the DeepLabv3+, however, the baseline performance is always higher for in-domain data (green) and fluctuates for OOD data (red), but ultimately the baseline performance is also higher for OOD data in the last iterations.

Since the mean cosine similarity did not provide any insights into the causes for the poor performance of the re-parameterized DeepLabv3+, we further investigated the mean cosine similarity for individual layers  $\ell$  of the networks as shown in Figure 6. We show the layer-wise mean cosine similarity for the encoder network ( $\theta_\ell^E$ ), where  $\ell$  indicates the layer index. For the purpose of clarity, we only mark the first layer of each of the major network blocks. For the SegFormer, we indicate the transformer blocks by  $Bb$  with  $b$  being the block index (cf. Fig-

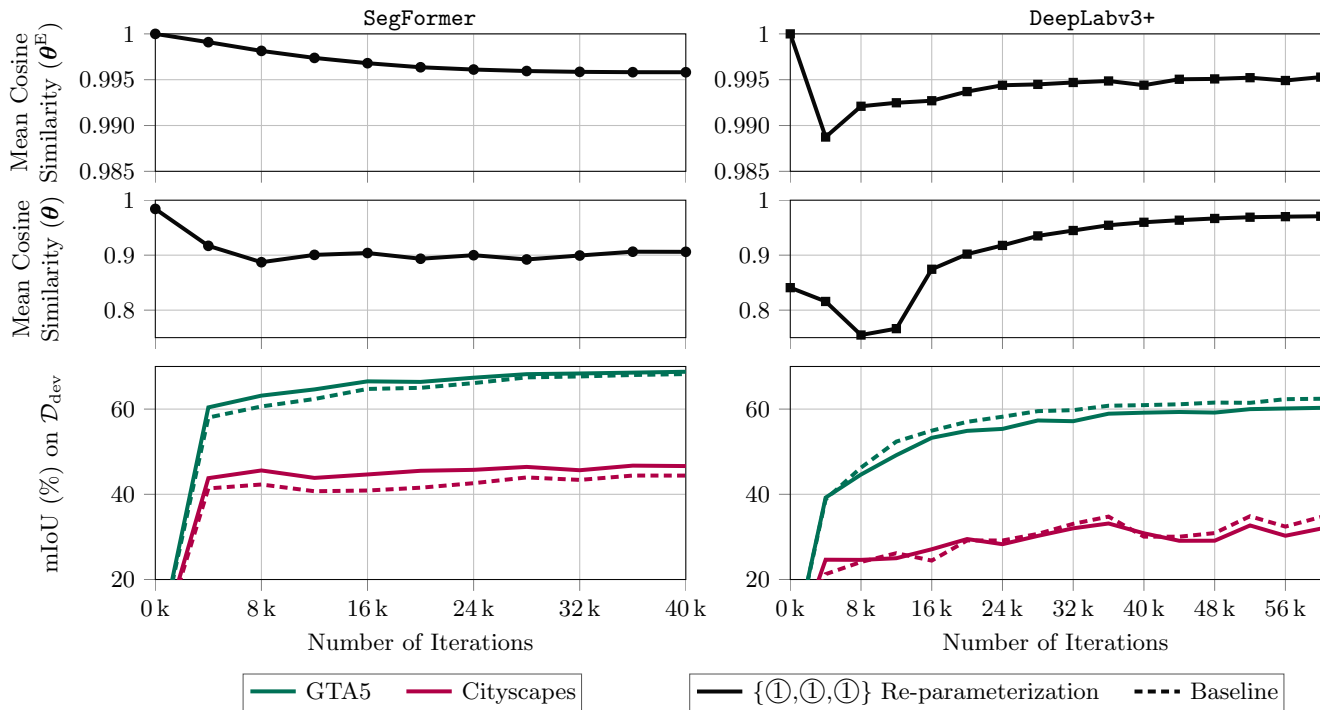


Figure 5. Comparison of **mean cosine similarity vs. mIoU performance** for SegFormer and DeepLabv3+ during training. The first row shows the mean cosine similarity for the encoder only ( $\theta^E$ ), the second row for the full network ( $\theta^E$ ). The mean cosine similarity is computed between three models. In the bottom row, the mIoU is given for the dev sets of GTA5 (green) and Cityscapes (red).

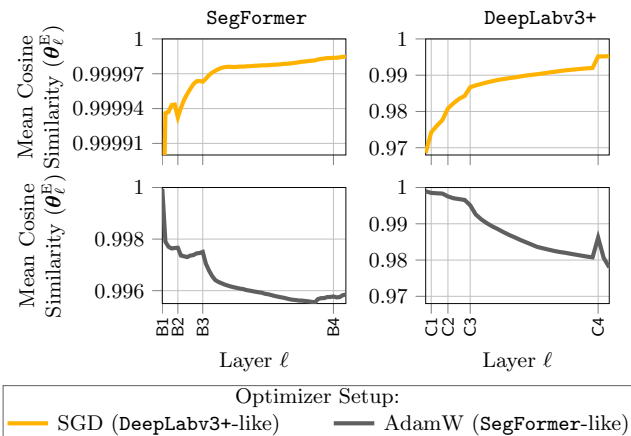


Figure 6. Comparison of the **layer-wise mean cosine similarity** for the encoder only ( $\theta_\ell^E$ ). Shown are results for the SegFormer (left) and DeepLabv3+ (right) architectures, which were trained with the standard DeepLabv3+ or SegFormer optimizer setup, shown in gray and yellow, respectively. The mean cosine similarity is computed between three models. For the purpose of clarity, we only indicate the individual network blocks ( $B_b$  for SegFormer,  $C_c$  for DeepLabv3+) on the x axis.

ure 8). For the DeepLabv3+, we indicate the convolutional blocks, as defined by He et al. [48], by  $C_c$  with  $c$

being the block index. *It can be seen that the choice of optimizer has a significant effect on the layer-wise cosine similarity.* For the standard DeepLabv3+ optimizer setup with SGD, the cosine similarity for both network architectures is higher in deeper layers and lower in earlier layers. In contrast, when the standard SegFormer optimizer setup with AdamW is employed, the cosine similarity is highest for earlier layers and drops for deeper layers. This specific property might be important for a well-performing encoder re-parameterization. As already shown in Table 4, this optimizer setup (AdamW [24]) also allows the DeepLabv3+ to improve over its baseline performance. Accordingly, we used the AdamW optimizer for our ReVT method in the main paper.

## D. Additional Ablation Studies

In this section, we will investigate the weighting of the base models and compare the ReVT re-parameterization w.r.t. re-parameterized network parts, layers, and the number of base models. Afterwards, we evaluate different base model augmentations and optimizer methods during training to design our final ReVT.

**Weighting of networks:** In Figure 7 we depict multiple possible weighting combinations for three models with the best combination (marked with a blue circle) achieving an mIoU of 47.68%, while the uniform re-parameterization

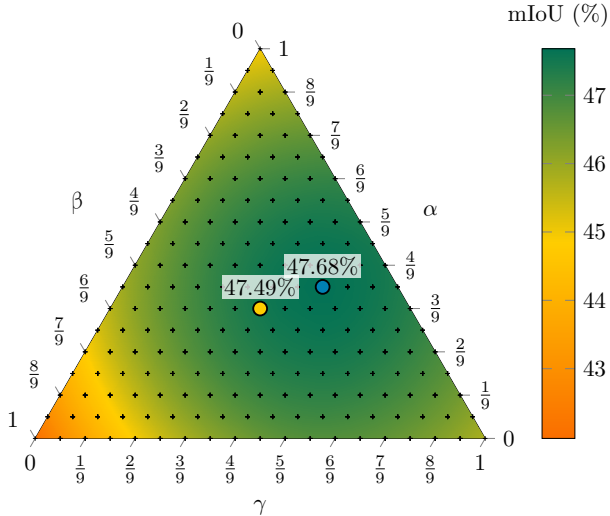


Figure 7. Ternary plot showing the performance (mIoU (%)) of three baseline models  $\{\textcircled{1}, \textcircled{1}, \textcircled{1}\}$  and multiple weight combinations  $(\alpha, \beta, \gamma)$ . The mIoU is calculated by a re-parameterization, where the parameters of the three models are weighted by the values of  $\alpha$ ,  $\beta$ , and  $\gamma$ , respectively. The **training** of the base models (SegFormer) was performed on the synthetic **GTA5** training set ( $\mathcal{D}^S = \mathcal{D}_{\text{train}}^{\text{GTA5}}$ ). The **evaluation** is performed on the **Cityscapes development set** ( $\mathcal{D}^T = \mathcal{D}_{\text{dev}}^{\text{CS}}$ ). The center of the plot ( $\alpha = \beta = \gamma = \frac{1}{3}$ ) corresponds to a uniform encoder re-parameterization.

(all models weighted with  $\frac{1}{3}$ , marked with a yellow circle) achieves an mIoU of 47.49%. It can be seen that the weighting of the individual models is actually quite insensitive, which is why we decide to use the simple variant of the uniform encoder re-parameterization (i.e., weights  $\frac{1}{3}$ ).

**ReVT vs. ensembles:** In Table 10 we compare the ReVT and a network ensemble for different combinations of training settings. It can be seen that the ReVT method outperforms the network ensemble not only for the combination of three baseline models  $\{\textcircled{1}, \textcircled{1}, \textcircled{1}\}$ , as already shown in Figure 4, but also for all other tested base model combinations  $\{\textcircled{u}_1, \textcircled{u}_2, \textcircled{u}_3\}$  on the real test\* sets. However, the network ensemble is slightly better in the synthetic source domain ( $\mathcal{D}_{\text{dev}}^{\text{GTA5}}$ ), which has little relevance for real-world applications. The same applies to the base model combination  $\{\textcircled{5}, \textcircled{5}, \textcircled{5}\}$ , where the ensemble performs slightly better on the SYNTHIA dataset ( $\mathcal{D}_{\text{dev}}^{\text{SYN}}$ ). The overall best performance for each dataset is highlighted in light green. It can be seen that the ReVT method achieves top performance for each dataset and is even on par with the ensemble on the source dataset ( $\mathcal{D}_{\text{dev}}^{\text{GTA5}}$ ). In general, once again, the ReVT  $\{\textcircled{1}, \textcircled{4}, \textcircled{6}\}$  seems to be a strong ReVT, generalizing well to unseen datasets.

Table 10. Performance (mIoU (%)) of a **network ensemble** vs. the **ReVT** for different base model combinations. The **training** of the base models (SegFormer) was performed on the **GTA5** ( $\mathcal{D}^S = \mathcal{D}_{\text{train}}^{\text{GTA5}}$ ) dataset. The **evaluation** is performed on the **GTA5** and **SYNTHIA development sets** and on the **test\* data** of various real-world target datasets ( $\mathcal{D}^T = \mathcal{D}_{\text{test}*}$ ). Reported is the mean value over the respective number of employed base models. Best result for each base model combination in bold face, overall best performance per dataset is highlighted in light green.

Base Models	Method	mIoU (%) on				
		$\mathcal{D}_{\text{dev}}^{\text{GTA5}}$	$\mathcal{D}_{\text{dev}}^{\text{SYN}}$	$\mathcal{D}_{\text{test}*}^{\text{CS}}$	$\mathcal{D}_{\text{test}*}^{\text{BDD}}$	$\mathcal{D}_{\text{test}*}^{\text{MV}}$
$\{\textcircled{1}, \textcircled{1}, \textcircled{1}\}$	Ensemble	<b>68.8</b>	34.6	46.9	44.8	48.0
	ReVT	68.6	<b>35.5</b>	<b>49.3</b>	<b>45.3</b>	<b>49.3</b>
$\{\textcircled{2}, \textcircled{2}, \textcircled{2}\}$	Ensemble	<b>69.2</b>	33.1	43.8	43.9	46.6
	ReVT	69.1	<b>34.1</b>	<b>44.9</b>	<b>44.1</b>	<b>47.6</b>
$\{\textcircled{3}, \textcircled{3}, \textcircled{3}\}$	Ensemble	<b>69.7</b>	33.8	43.8	42.6	47.3
	ReVT	<b>69.7</b>	<b>35.0</b>	<b>45.7</b>	<b>43.5</b>	<b>48.8</b>
$\{\textcircled{4}, \textcircled{4}, \textcircled{4}\}$	Ensemble	<b>65.9</b>	<b>36.8</b>	47.5	47.4	52.6
	ReVT	65.7	36.2	<b>48.8</b>	<b>47.5</b>	<b>53.2</b>
$\{\textcircled{5}, \textcircled{5}, \textcircled{5}\}$	Ensemble	<b>68.6</b>	35.1	47.5	45.4	50.4
	ReVT	<b>68.5</b>	<b>35.9</b>	<b>48.6</b>	<b>45.9</b>	<b>51.2</b>
$\{\textcircled{6}, \textcircled{6}, \textcircled{6}\}$	Ensemble	<b>65.0</b>	<b>37.0</b>	48.0	47.9	52.8
	ReVT	64.9	36.1	<b>49.7</b>	<b>48.5</b>	<b>53.5</b>
$\{\textcircled{4}, \textcircled{5}, \textcircled{6}\}$	Ensemble	<b>67.2</b>	35.1	48.1	47.3	51.6
	ReVT	<b>66.4</b>	<b>36.9</b>	<b>49.5</b>	<b>48.1</b>	<b>53.1</b>
$\{\textcircled{1}, \textcircled{4}, \textcircled{6}\}$	Ensemble	<b>68.5</b>	34.2	47.0	45.9	50.1
	ReVT	66.4	<b>37.3</b>	<b>50.0</b>	<b>48.0</b>	<b>52.8</b>

## E. ReVT with SYNTHIA as Source

In this section, we provide additional results for our method, when trained with the SYNTHIA dataset as source domain:  $\mathcal{D}^S = \mathcal{D}_{\text{train}}^{\text{SYN}}$ . In Table 11 we evaluate the various augmentation methods  $\textcircled{a}$  that were already evaluated for models trained on GTA5 ( $\mathcal{D}^S = \mathcal{D}_{\text{train}}^{\text{SYN}}$ ) in Section 5. In the lower part of the table we report some ( $M = 3$ ) combinations  $\{\textcircled{u}_1, \textcircled{u}_2, \textcircled{u}_3\}$  of these base models by our re-parameterization. The gray columns indicate our development sets ( $\mathcal{D}_{\text{dev}}$ ), where the light gray column is  $\mathcal{D}_{\text{dev}}^{\text{SYN}}$ . The OOD mean mIoU of  $\mathcal{D}_{\text{dev}}^{\text{GTA5}}$  and  $\mathcal{D}_{\text{dev}}^{\text{CS}}$  is shown in the dark gray columns. It can be seen in the upper part of Table 11 that the augmentation methods do not improve the performance as much as for models trained on GTA5. The best OOD mean performance is achieved with the baseline model  $\textcircled{1}$ . On the test\* mean the PixMix\* augmentation works best, followed by the combination of PixMix\* and the bilateral filter, and baseline model.

Although the individual augmentations do not perform as well for these models, the ReVT  $\{\textcircled{1}, \textcircled{4}, \textcircled{6}\}$ , which we

Table 11. Performance (mIoU (%)) of the SegFormer model (with an MiT-B5 encoder) using different domain generalization methods. **Training** was performed on the synthetic SYNTHIA ( $\mathcal{D}^S = \mathcal{D}_{\text{train}}^{\text{SYN}}$ ) dataset. **Evaluation** is performed on the **Cityscapes**, **GTA5**, and **SYNTHIA development sets** (gray columns) and on the **test\* data** of various real-world target datasets ( $\mathcal{D}^T = \mathcal{D}_{\text{test}^*}$ ). Reported is the mean mIoU  $\pm$  the standard deviation of  $M=3$  models with various image augmentations. For the ReVT, the mean  $\pm$  standard deviation is computed with one averaged encoder and the three associated decoders  $m \in \{1, 2, 3\}$ . For models trained on SYNTHIA, we evaluate over 16 classes, as is common practice [19]. Best results in bold face, second-best underlined.

Method performed:		mIoU (%) on							
		$\mathcal{D}_{\text{dev}}^{\text{SYN}}$	$\mathcal{D}_{\text{dev}}^{\text{GTA5}}$	$\mathcal{D}_{\text{dev}}^{\text{CS}}$	<b>OOD mean</b>	$\mathcal{D}_{\text{test}^*}^{\text{CS}}$	$\mathcal{D}_{\text{test}^*}^{\text{BDD}}$	$\mathcal{D}_{\text{test}^*}^{\text{MV}}$	<b>test* mean</b>
...during training	Baseline ①	76.5 $\pm$ 0.1	42.8 $\pm$ 0.4	<b>44.3</b> $\pm$ 1.3	<b>43.5</b> $\pm$ 1.2	<b>45.1</b> $\pm$ 1.6	35.2 $\pm$ 1.4	<b>42.5</b> $\pm$ 0.9	<u>40.9</u> $\pm$ 4.4
	-PhotoAug ②	<u>77.3</u> $\pm$ 0.1	39.8 $\pm$ 0.6	41.4 $\pm$ 0.5	40.6 $\pm$ 0.9	41.6 $\pm$ 0.8	33.7 $\pm$ 1.0	40.8 $\pm$ 0.6	38.7 $\pm$ 3.6
	-PhotoAug, -Rand. Flip ③	<b>78.3</b> $\pm$ 0.0	40.7 $\pm$ 0.7	41.3 $\pm$ 0.3	41.0 $\pm$ 0.6	41.8 $\pm$ 0.2	34.3 $\pm$ 1.1	41.3 $\pm$ 0.3	39.1 $\pm$ 3.5
	+PixMix* [13] ④	73.8 $\pm$ 0.1	<b>43.1</b> $\pm$ 1.0	42.6 $\pm$ 1.7	42.8 $\pm$ 1.4	42.6 $\pm$ 2.2	<b>38.4</b> $\pm$ 0.8	<b>42.5</b> $\pm$ 0.6	<b>41.2</b> $\pm$ 2.4
	+Bilateral Filter (BF) [39] ⑤	76.0 $\pm$ 0.1	42.9 $\pm$ 0.9	43.2 $\pm$ 0.7	<u>43.1</u> $\pm$ 0.8	43.1 $\pm$ 0.9	36.0 $\pm$ 0.7	41.7 $\pm$ 0.4	40.3 $\pm$ 3.1
	+PixMix* [13] +BF [39] ⑥	72.3 $\pm$ 0.1	41.2 $\pm$ 0.6	<u>43.3</u> $\pm$ 0.4	42.2 $\pm$ 1.2	<u>43.0</u> $\pm$ 0.5	<u>38.0</u> $\pm$ 0.6	41.6 $\pm$ 0.3	<u>40.9</u> $\pm$ 2.2
...after training	ReVT {①,①,①}	76.2 $\pm$ 0.1	42.3 $\pm$ 0.1	44.9 $\pm$ 0.4	43.6 $\pm$ 1.3	<u>45.8</u> $\pm$ 0.4	35.8 $\pm$ 0.2	43.6 $\pm$ 0.2	41.7 $\pm$ 4.3
	ReVT {②,②,②}	<u>76.9</u> $\pm$ 0.0	40.4 $\pm$ 0.2	42.2 $\pm$ 0.2	41.3 $\pm$ 1.0	42.6 $\pm$ 0.1	34.7 $\pm$ 0.3	42.0 $\pm$ 0.3	39.8 $\pm$ 3.6
	ReVT {③,③,③}	<b>78.0</b> $\pm$ 0.1	41.3 $\pm$ 0.6	41.9 $\pm$ 0.3	41.6 $\pm$ 0.6	42.5 $\pm$ 0.4	35.2 $\pm$ 0.3	42.7 $\pm$ 0.3	40.1 $\pm$ 3.5
	ReVT {④,④,④}	73.7 $\pm$ 0.1	<u>43.4</u> $\pm$ 0.1	44.1 $\pm$ 0.3	43.7 $\pm$ 0.4	44.3 $\pm$ 0.3	39.5 $\pm$ 0.3	<u>44.0</u> $\pm$ 0.3	42.6 $\pm$ 2.2
	ReVT {⑤,⑤,⑤}	75.7 $\pm$ 0.0	<b>44.1</b> $\pm$ 0.7	44.4 $\pm$ 0.3	<b>44.2</b> $\pm$ 0.6	44.5 $\pm$ 0.3	37.1 $\pm$ 0.3	43.0 $\pm$ 0.2	41.5 $\pm$ 3.2
	ReVT {⑥,⑥,⑥}	72.2 $\pm$ 0.1	42.4 $\pm$ 0.2	44.4 $\pm$ 0.2	43.4 $\pm$ 1.0	44.3 $\pm$ 0.2	38.9 $\pm$ 0.2	42.9 $\pm$ 0.0	42.0 $\pm$ 2.3
	ReVT {④,⑤,⑥}	74.0 $\pm$ 0.6	43.2 $\pm$ 0.4	<u>45.0</u> $\pm$ 0.4	<u>44.1</u> $\pm$ 0.9	45.1 $\pm$ 0.4	<u>39.6</u> $\pm$ 0.4	<u>44.0</u> $\pm$ 0.2	<u>42.9</u> $\pm$ 2.4
	ReVT {①,④,⑥}	74.1 $\pm$ 0.6	42.7 $\pm$ 0.3	<b>45.7</b> $\pm$ 0.5	<b>44.2</b> $\pm$ 1.5	<b>46.3</b> $\pm$ 0.3	<b>40.3</b> $\pm$ 0.5	<b>44.8</b> $\pm$ 0.1	<b>43.8</b> $\pm$ 2.6

already identified as our best ReVT in Section 5, provides both top OOD mean (44.2%) and test\* mean (43.8%) performance. Again, second-best results are achieved with the ReVT {④,⑤,⑥}. Additionally, the ReVT {⑤,⑤,⑤} is on par with the ReVT {①,④,⑥} for the OOD mean performance when trained on SYNTHIA.

In the following, we compare also against prior art that have also been evaluated with SYNTHIA as source domain. Again, we choose the ReVT {①,④,⑥} and ReVT {④,⑤,⑥} for the comparison with prior art. The results are shown in Table 12. In contrast to the models trained on GTA5, for models trained on SYNTHIA, the ReVT {①,④,⑥} does not always reach the top BM mean performance. Similar to the GTA5-trained models, the performance of both ReVT variants remains slightly behind that of the baseline for synthetic source domain data ( $\mathcal{D}_{\text{dev}}^{\text{SYN}}$ ), which, however, has little relevance for practical real-world applications.

For the small (group 1) and mid-sized (group 2) models, the ReVT {④,⑤,⑥} yields a slightly better performance of 45.79% vs. 45.44% (baseline: 44.09%) and 48.99% vs. 48.84% (baseline: 46.72%), respectively. For the large models (group 3), the ReVT {①,④,⑥} yields the best performance with a BM mIoU of 49.64% (base-

line: 48.42%). In summary, on the benchmark (BM) data, our proposed SYNTHIA-trained ReVT models achieve an mIoU improvement of +1.2% absolute (large models) to +1.7% absolute (small models).

It should be noted that no prior work reported on all datasets necessary for the benchmark (BM) mean when trained with SYNTHIA as source domain. All of our ReVTs improve on the prior art for the reported domains. Only for the smallest models in group 1 the SAN+SAW method [30] achieves a higher mIoU on the BDD dataset (best prior art: 35.42% vs. ours: 35.18%). For the mid-sized models we already improve on this domain (best prior art: 37.40% vs. ours: 38.73%), and interestingly significantly excel the SAN+SAW method (ours: 38.73% vs. SAN+SAW: 35.98%).

## F. SegFormer Block Diagrams

In Section 5 we investigated the effect of the re-parameterization on different network parts (cf. Table 3) and block or layer types (cf. Table 5). To give the reader a better idea of how the network is structured and where the individual block and layer types are located in the network, an hierarchically illustrated overview of the SegFormer architecture with an MiTB5 encoder is given in Figures 8, 9,

Table 12. Performance (mIoU (%)) of various domain generalization methods employing different segmentation networks, sorted into three performance groups. **Training** was performed on the synthetic **SYNTHIA** ( $\mathcal{D}^S = \mathcal{D}_{\text{train}}^{\text{SYN}}$ ) dataset. The results marked with  $^\circ$  are cited from [20] and with  $^*$  are cited from the respective paper. All results without any identifier are simulated. **Evaluation** is performed on the **SYNTHIA** and **GTA5 development sets** and on the **test\* data** of various real-world target datasets ( $\mathcal{D}^T = \mathcal{D}_{\text{test}^*}$ ). BM means benchmark. For our simulations we report mean values over three runs with different seeding. For models trained on SYNTHIA, we evaluate over 16 classes, as is common practice [19]. Best performance per group in bold face, second best underlined.

Enc.	Method	$ \theta $ ( $\cdot 10^6$ )	Single Source	Frame Rate [fps]	mIoU (%) on								
					$\mathcal{D}_{\text{test}^*}^{\text{CS}}$	$\mathcal{D}_{\text{test}^*}^{\text{BDD}}$	$\mathcal{D}_{\text{test}^*}^{\text{MV}}$	$\mathcal{D}_{\text{dev}}^{\text{SYN}}$	$\mathcal{D}_{\text{dev}}^{\text{GTA5}}$	$\mathcal{D}_{\text{test}^*}^{\text{ACDC}}$	$\mathcal{D}_{\text{test}^*}^{\text{KIT}}$	BM mean	
Group 1	ResNet-50	Baseline $^\circ$	49.6	✓	7.9	28.36	25.16	27.24	-	-	-	-	-
		DRPC $^\circ$ [44]	49.6	✗	8.3	35.65	31.53	32.74	-	-	-	-	-
		SAN+SAW $^*$ [30]	25.6	✓	8.1	38.92	<b>35.42</b>	34.52	-	29.16	-	-	-
	MiT-B2	Baseline	27.4	✓	<b>12.0</b>	39.71	29.76	38.37	<b>74.78</b>	37.83	26.16	<u>35.18</u>	44.09
		Ours: ReVT {④,⑤,⑥}	27.4	✓	<b>12.0</b>	<b>41.09</b>	<u>35.18</u>	<u>40.21</u>	<u>71.59</u>	<b>40.88</b>	<b>30.39</b>	34.64	<b>45.79</b>
		Ours: ReVT {①,④,⑥}	27.4	✓	<b>12.0</b>	<u>40.91</u>	34.53	<b>40.44</b>	71.45	<u>39.87</u>	<u>30.13</u>	<b>35.29</b>	<u>45.44</u>
Group 2	ResNet-101	Baseline $^\circ$	68.6	✓	7.9	29.67	25.64	28.73	-	-	-	-	-
		DRPC $^\circ$ [44]	68.6	✗	5.3	37.58	34.34	34.12	-	-	-	-	-
		FSDR $^*$ [14]	68.6	✗	5.3	40.80	37.40	39.60	-	-	-	-	-
		SAN+SAW $^*$ [30]	44.6	✓	5.3	40.87	35.98	37.26	-	30.79	-	-	-
	MiT-B3	Baseline	47.2	✓	<b>10.7</b>	42.43	33.33	40.47	<b>75.82</b>	41.53	29.73	35.91	46.72
		Ours: ReVT {④,⑤,⑥}	47.2	✓	<b>10.7</b>	<b>45.26</b>	<b>38.73</b>	<u>42.86</u>	73.12	<b>44.99</b>	<b>35.27</b>	<b>36.42</b>	<b>48.99</b>
Ours: ReVT {①,④,⑥}		47.2	✓	<b>10.7</b>	<u>44.97</u>	<u>38.65</u>	<b>43.00</b>	<u>73.16</u>	<u>44.42</u>	<u>35.16</u>	<u>36.31</u>	<u>48.84</u>	
Group 3	MiT-B5	Baseline	84.7	✓	<b>9.7</b>	45.07	35.19	42.51	<b>76.49</b>	<u>42.82</u>	30.81	37.02	48.42
		Ours: ReVT {④,⑤,⑥}	84.7	✓	<b>9.7</b>	<u>45.08</u>	<u>39.62</u>	<u>43.99</u>	73.96	<b>43.25</b>	<u>35.12</u>	<u>37.20</u>	<u>49.18</u>
		Ours: ReVT {①,④,⑥}	84.7	✓	<b>9.7</b>	<b>46.28</b>	<b>40.30</b>	<b>44.76</b>	<u>74.11</u>	42.74	<b>35.75</b>	<b>37.86</b>	<b>49.64</b>

10, 11, and 12.

## G. Discussion of Limitations

Although modern methods for domain generalization provide good performance on completely unseen real data (after training on synthetic data), the performance still remains behind that of modern methods for unsupervised domain adaptation (UDA) [47,49]. Such a comparison, however, is not entirely fair, since UDA methods employ unlabeled data from a target domain (typically Cityscapes) during the training process, which we intentionally avoid in domain generalization. Nevertheless, it should be noted that better performance on a specific target domain can be achieved, if samples from this domain are available during training.

Our proposed method cannot be applied advantageously to any already trained model, since the optimizer choice has a significant impact on the performance. To be fair, however, this is the case with all prior art methods as well. Most of them additionally extend the training process considerably, far beyond the choice of the optimizer [5, 14, 20, 30, 44].

## H. Discussion of Ethical Implications

Although well generalizing semantic segmentation has many civilian applications that provide great value to society, e.g., automated driving, robotics, and medical applications, this technology can also be used for military and surveillance applications. Research on better generalizing methods may also indirectly contribute to the improvement of these applications.

Another aspect to consider are biases in the employed datasets. Three of the five real datasets (Cityscapes, ACDC, KITTI) were captured in Central Europe, one in the USA (BDD100k), and only one contains data from all over the world (Mapillary Vistas). This may lead to biases regarding different ethnicities in the data, which were not investigated further in this paper. For the reported results on improved generalization from synthetic to real data, the biases may be negligible, but should be considered for possible real-world applications.

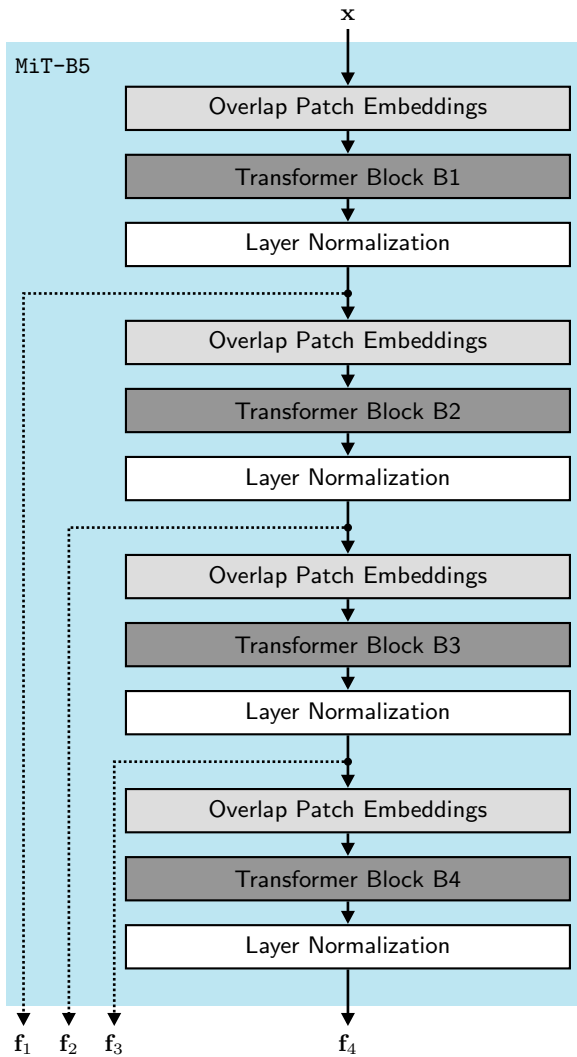


Figure 8. Overview of the MiT-B5 encoder that is employed by the largest SegFormer model. This is the standard encoder employed in the segmentation model (cf. Figure 2).

### Additional References

- [47] Nikita Araslanov and Stefan Roth. Self-Supervised Augmentation Consistency for Adapting Semantic Segmentation. In *Proc. of CVPR*, pages 15384–15394, virtual, June 2021.
- [48] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proc. of CVPR*, pages 770–778, Las Vegas, NV, USA, June 2016
- [49] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. DAFormer: Improving Network Architectures and Training Strategies for Domain-Adaptive Semantic Segmentation. In *Proc. of CVPR*, pages 9924–9935, New Orleans, LA, USA, June 2022

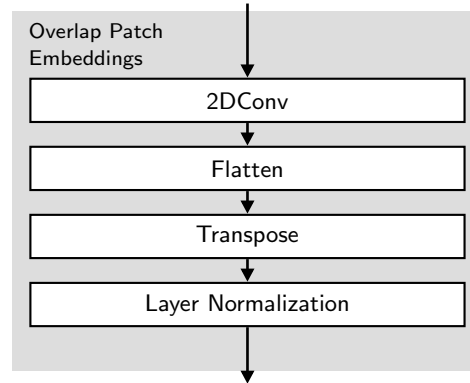


Figure 9. Overview of the overlap patch embeddings block that is employed in the MiT encoder (cf. Figure 8 for MiT-B5).

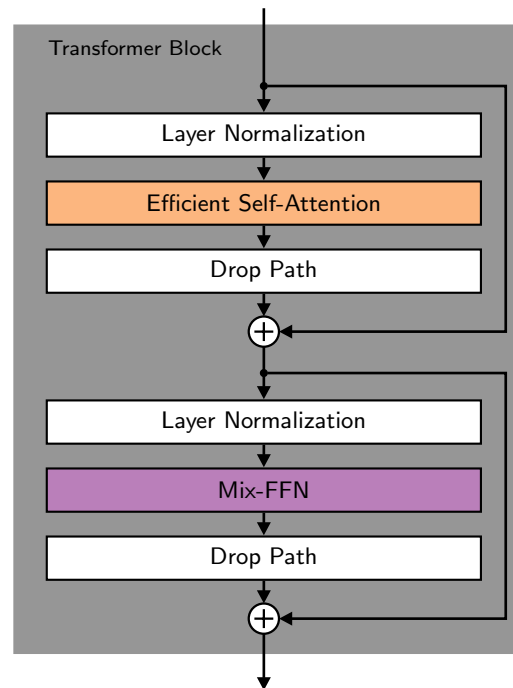


Figure 10. Overview of the transformer block that is employed in the MiT encoder (cf. Figure 8 for MiT-B5).

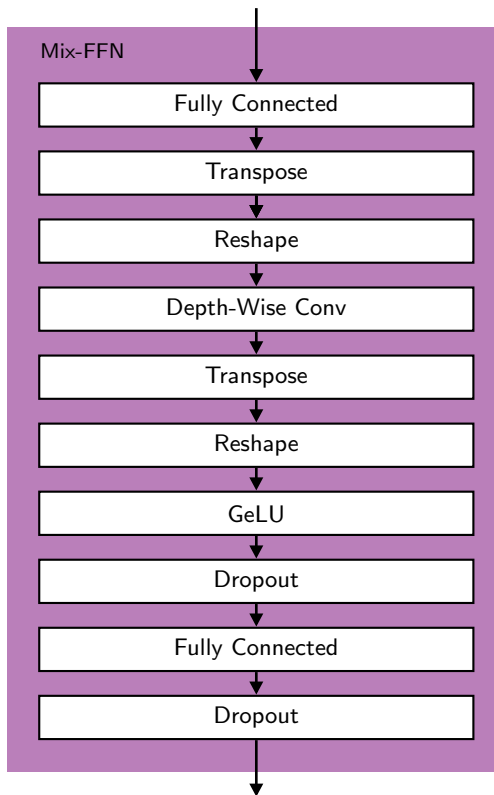


Figure 11. Overview of the Mix-FFN block that is employed in the transformer block (cf. Figure 10).

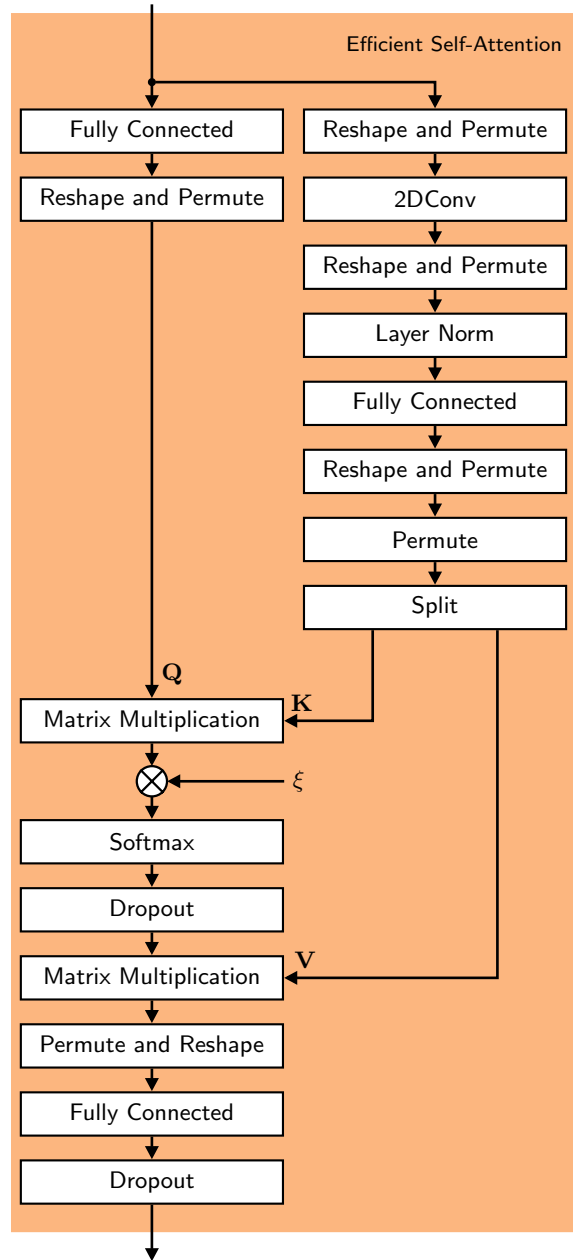


Figure 12. Overview of the efficient self-attention block that is employed in the transformer block (cf. Figure 10). The fixed scaling factor  $\xi$  is a hyperparameter and block-dependent.