# Masked Discriminators for Content-Consistent Unpaired Image-to-Image Translation

Bonifaz Stuhr, Jürgen Brauer, Bernhard Schick, and Jordi Gonzàlez



PFD→Cityscapes

Day→Night

Viper→Cityscapes

Clear→Snowy

Fig. 1: Results of our method.

*Abstract*—A common goal of unpaired image-to-image translation is to preserve content consistency between source images and translated images while mimicking the style of the target domain. Due to biases between the datasets of both domains, many methods suffer from inconsistencies caused by the translation process. Most approaches introduced to mitigate these inconsistencies do not constrain the discriminator, leading to an even more ill-posed training setup. Moreover, none of these approaches is designed for larger crop sizes. In this work, we show that masking the inputs of a global discriminator for both domains with a content-based mask is sufficient to reduce content inconsistencies significantly. However, this strategy leads to artifacts that can be traced back to the masking process. To reduce these artifacts, we introduce a local discriminator that operates on pairs of small crops selected with a similarity sampling strategy. Furthermore, we apply this sampling strategy to sample global input crops from the source and target dataset. In addition, we propose feature-attentive denormalization to selectively incorporate content-based statistics into the generator stream. In our experiments, we show that our method achieves state-of-the-art performance in photorealistic sim-to-real translation and weather translation and also performs well in day-to-night translation. Additionally, we propose the cKVD metric, which builds on the sKVD metric and enables the examination of translation quality at the class or category level.

*Index Terms*—masked discriminators, feature-attentive denormalization, generative adversarial networks (GANs), content-consistent, unpaired image-to-image translation

Bonifaz Stuhr and Jordi Gonzàlez are with the Department of Computer Science, Autonomous University of Barcelona, 08193 Bellaterra, Barcelona, Spain (e-mail: bonifaz.stuhr@hs-kempten.de; jordi.gonzalez@uab.cat).

Jürgen Brauer is with the Department of Computer Science, University of Applied Sciences Kempten, 87435 Kempten, Germany (e-mail: juergen.brauer@hs-kempten.de).

Bernhard Schick and Bonifaz Stuhr are with the IFM, Junkersstraße 1A, 87734 Benningen, Germany (e-mail: bernhard.schick@hs-kempten.de).

## I. Introduction

UNPAIRED image-to-image translation aims at transferring images from a source domain to a target domain when no paired examples are given. Recently, this field has attracted increasing interest and has advanced several use cases, such as photorealism [1]–[3], neural rendering [4], domain adaptation [5], [6], the translation of seasons or daytime [1], [7], [8], and artistic style transfer [9]–[11]. Current work has primarily focused on improving translation quality [8], [12], efficiency [13], [14], multi-modality [15], [16], and content consistency [2], [3]. Due to the ill-posed nature of the unpaired image-to-image translation task and biases between datasets, content consistency is difficult to achieve. To mitigate content inconsistencies, several methods have been proposed that constrain the generator of GANs [15]–[24]. However, only constraining the generator leads to an unfair setup, as biases in the datasets can be detected by the discriminator: The generator tries to achieve content consistency by avoiding biases in the output, while the discriminator is still able to detect biases between both datasets and, therefore, forces the generator to include these biases in the output, for example, through hallucinations. Constraining the discriminator [2], [25], [26] or improving the sampling of training pairs [2], [27] is currently underexplored, especially for content consistency on a global level, where the discriminator has a global view on larger image crops instead of a local view on small crops. In this work, we propose *masked conditional discriminators*, which operate on masked global crops of the inputs to mitigate content inconsistencies. We combine these discriminators with an efficient sampling strategy based on a pre-trained robust segmentation model to sample similar global crops. Furthermore, we argue that when transferring feature statistics from the content stream of the source image to the

generator stream, content-unrelated feature statistics from the content stream could affect image quality if the generator is unable to ignore this information since the output image should mimic the target domain. Therefore, we propose a *feature-attentive denormalization (FATE)* block that extends feature-adaptive denormalization (FADE) [7] with an attention mechanism. This block allows the generator to selectively incorporate statistical features from the content stream into the generator stream. In our experiments, we find that our method achieves state-of-the-art performance on most of the benchmarks shown in Figure 1.

Our contributions can be summarized as follows:

- We propose an efficient sampling strategy that utilizes robust semantic segmentations to sample similar global crops. This reduces biases between both datasets induced by semantic class misalignment.
- We combine this strategy with masked conditional discriminators to achieve content consistency while maintaining a more global field of view.
- We extend our method with an unmasked local discriminator. This discriminator operates on local, partially class-aligned patches to minimize the underrepresentation of frequently masked classes and associated artifacts.
- We propose a feature-attentive denormalization (FATE) block, which selectively fuses statistical features from the content stream into the generator stream.
- We propose the class-specific Kernel VGG Distance (cKVD) that builds upon the semantically aligned Kernel VGG Distance (sKVD) [2] and uses robust segmentations to incorporate class-specific content inconsistencies in the perceptual image quality measurement.
- In our experiments, we show that our method achieves state-of-the-art performance on photo-realistic sim-to-real transfer and the translation of weather and performs well for daytime translation.

## II. RELATED WORK

**Unpaired image-to-image translation.** Following the success of GANs [28], the conditional GAN framework [29] enables image generation based on an input condition. Pix2Pix [30] uses images from a source domain as a condition for the generator and discriminator to translate them to a target domain. Since Pix2Pix relies on a regression loss between generated and target images, translation can only be performed between domains where paired images are available. To achieve unpaired image-to-image translation, methods like CycleGAN [17], UNIT [22], and MUNIT [15] utilize a second GAN to perform the translation in the opposite direction and impose a cycle-consistency constraint or weight-sharing constraint between both GANs. However, these methods require additional parameters for the second GAN, which are used to learn the unpaired translation and are omitted when inferring a one-sided translation. In works such as TSIT [7] and CUT [31], these additional parameters are completely omitted at training time by either utilizing a perceptual loss [32] between the input image of the generator

and the image to be translated or by patchwise contrastive learning. Recently, additional techniques have achieved promising results, like pseudo-labeling [4] or a conditional discriminator based on segmentations created with a robust segmentation model for both domains [2]. Furthermore, there are recent efforts to adapt diffusion models to unpaired image-to-image translation [33]–[35].

**Content consistency in unpaired image-to-image translation.** Due to biases between unpaired datasets, the content of translated samples can not be trivially preserved [2]. There are ongoing efforts to preserve the content of an image when it is translated to another domain by improving various parts of the training pipeline: Several consistency constraints have been proposed for the generator, which operate directly on the translated image [16]–[18], on a transformation of the translated image [19], [20], [24], [36], [37], or on distributions of multi-modal translated images [21]. The use of a perceptual loss [32] or LPIPS loss [38] between input images and translated images, as in [7] and [2], can also be considered a consistency constraint between transformed images. In [39] content consistency is enforced with self-supervised in-domain and cross-domain patch position prediction. There are works that enforce consistency by constraining the latent space of the generator [15], [22], [23]. Semantic scene inconsistencies can be mitigated with a separate segmentation model [16], [24]. To avoid inconsistency arising from style transfer, features from the generator stream are masked before AdaIN [9], [40]. Another work exploits small perturbations in the input feature space to improve semantic robustness [3]. However, if the datasets of both domains are unbalanced, discriminators can use dataset biases as learning shortcuts, which leads to content inconsistencies. Therefore, only constraining the generator for content consistency still results in an ill-posed unpaired image-to-image translation setup. Constraining discriminators to achieve content consistency is currently underexplored, but recent work has proposed promising directions. There are semantic-aware discriminator architectures [2], [4], [25], [41] that enforce discriminators to base their predictions on semantic classes, or VGG discriminators [2], which additionally operate on abstract features of a frozen VGG model instead of the input images. Training discriminators with small patches [2] is another way to improve content consistency. To mitigate dataset biases during training for the whole model, sampling strategies can be applied to sample similar patches from both domains [2], [27]. Furthermore, in [26], a model is trained to generate a hyper-vector mapping between source and target images with an adversarial loss and a cyclic loss for content consistency. In contrast, our work utilizes a robust semantic mask to mask global discriminators with a large field of view, which provide the generator with the gradients of the unmasked regions. This leads to a content-consistent translation while preserving the global context. We combine this discriminator with an efficient sampling method that uses robust semantic segmentations to sample similar crops from both domains.
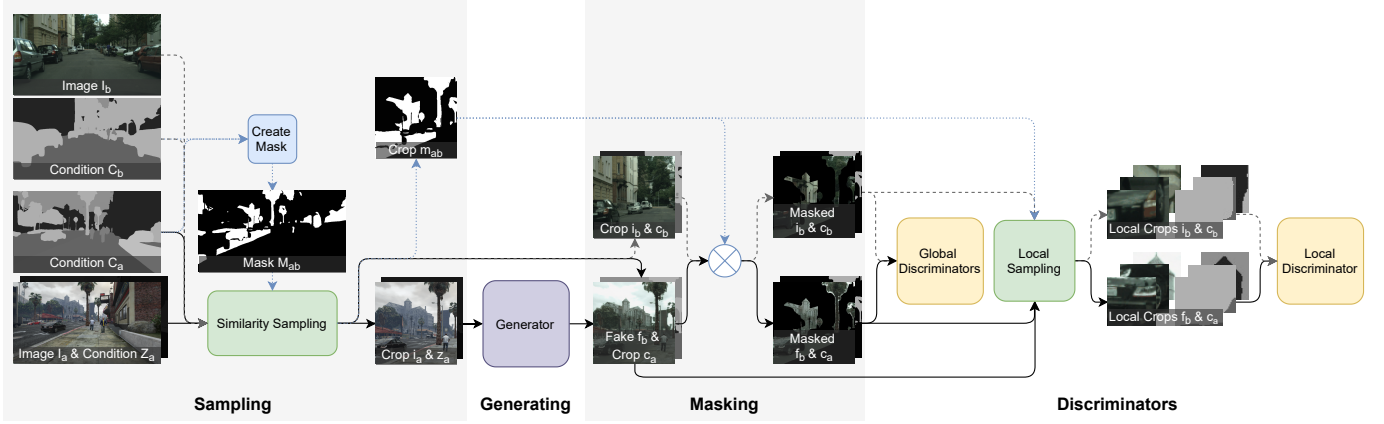
Fig. 2: **Method overview.** In our method, similar image crops from both domains ($i_a$, $i_b$) and their corresponding conditions ($c_a$, $c_b$, $z_a$) are selected via a sampling procedure. In this sampling procedure, a mask $M_{ab}$ is created from the conditions $C_a$ and $C_b$. This mask is used to sample crops from both datasets for which the semantic classes align by at least 50%. The cropped mask $m_{ab}$ is also used to mask the generated fake image $f_b$, the real images $i_b$, and the corresponding conditions for the global conditional discriminators. Through the mask, these discriminators can only see the parts of the crop where the semantic classes align. To further improve image quality, a local discriminator is introduced that works on a batch of small patches selected from the crop using our sampling technique. This discriminator is not masked and works on patches where the semantic classes do not fully align.

**Attention in image-to-image translation.** Previous work has utilized attention for different parts of the GAN framework. A common technique is to create attention mechanisms that allow the generator or discriminator to focus on important regions of the input [11], [42]–[45] or to capture the relationship between regions of the input(s) [10], [46], [47]. Other works guide a pixel loss with uncertainty maps computed from attention maps [48], exploit correlations between channel maps with scale-wise channel attention [46], disentangle content and style with diagonal attention [49], or merge features from multiple sources with an attentional block before integrating them into the generator stream [50]. In [51], an attention-based discriminator is introduced to guide the training of the generator with attention maps. Furthermore, ViTs [52] are adapted for unpaired image-to-image translation [53], [54], and the computational complexity of their self-attention mechanism is reduced for high-resolution translation [54]. In contrast, our work proposes an attention mechanism to selectively integrate statistics from the content stream of the source image into the generator stream. This allows the model to focus on statistical features from the content stream that are useful for the target domain.

## III. METHOD

We propose an end-to-end framework for unpaired image-to-image translation that transfers an image $I_a \in \mathbb{R}^{3 \times h \times w}$ from a source domain $a$ to an image $F_b \in \mathbb{R}^{3 \times h \times w}$ from a target domain $b$. Our goal is to design a method for content-consistent translations that utilizes a simple masking strategy for the global crops seen by the discriminators. We achieve this by combining an efficient segmentation-based sampling method that samples large crops from the input image with a masked discriminator that operates on these global crops. This is in contrast to EPE [2], which achieves content-consistent

translation at the local level by sampling small, similar image crops from both domains. To further improve image quality, we use a local discriminator that operates on a batch of small image patches sampled from the global input crops utilizing our sampling method. An overview of our method is shown in Figure 2. Furthermore, we propose a feature-attentive denormalization (FATE) block that extends feature-adaptive denormalization (FADE) [7] with an attention mechanism, allowing the generator to selectively incorporate statistical features from the content stream of the source image into the generator stream.

### A. Contend-based Similarity Sampling

To minimize the bias between both datasets in the early stage of our method, we sample similar image crops with an efficient sampling procedure. This procedure uses the one-hot encoded semantic segmentations $C_a \in \mathbb{R}^{d \times h \times w}$ and $C_b \in \mathbb{R}^{d \times h \times w}$ of both domains, where $d$ is the channel dimension of the one-hot encoding. In our case, these segmentations are created with the robust pre-trained MSeg model [55]. First, a mask $M_{ab} \in \mathbb{R}^{1 \times h \times w}$ is computed from the segmentations:

$$M_{ab} = \max_d(C_a \circ C_b), \qquad (1)$$

where $\circ$ denotes the Hadamard product. We can now sample semantically aligned image crops $i_a$ and $i_b$ from the images $I_a$ and $I_b$ with the crop $m_{ab}$ from mask $M_{ab}$. Thereby, we calculate the percentage of overlap of semantic classes between both image crops as follows:

$$\mathcal{P}_{match}(i_a) = \{i_b \mid \text{mean}(m_{ab}) > t\}, \qquad (2)$$

where $t$ is the similarity sampling threshold. In our case, we sample crops where more than 50% of the semantic classes align ($t > 0.5$). We use this procedure to sample crops $c_a$,

$c_b$, and $z_b$ from the discriminator conditions $C_a$, $C_b$, and the generator condition $Z_b$ as well. The cropped mask $m_{ab}$ is also used for our masked conditional discriminator.

### B. Contend-based Discriminator Masking

To train a discriminator with a global field of view that facilitates the usage of global properties of the scene, while simultaneously maintaining content consistency, we mask the discriminator input from both domains with a content-based mask $m_{ab}$. This mask erases all pixels from the discriminator input where the semantic classes do not align. This removes the bias between both datasets caused by the underlying semantic class distribution of the two domains without directly restricting the generator. The objective function of a conditional GAN with a masked discriminator that transfers image crops $i_a$ to domain $b$ can be then defined as follows:

$$\mathcal{L}_{madv} = \; \mathbb{E}_{i_b,c_b,m_{ab}}[\log D(i_b \circ m_{ab}|c_b \circ m_{ab})]$$
$$+ \; \mathbb{E}_{i_a,z_a,c_a,m_{ab}}[\log(1 - D(G(i_a|z_a) \circ m_{ab}|c_a \circ m_{ab}))]. \quad (3)$$

To ensure that the discriminator does not use the segmentation maps as learning shortcuts, we follow [2] and create the segmentations of both datasets using a robust segmentation model such as MSeg [55]. With this setting, we are able to train discriminators with large crop sizes with significantly reduced hallucinations in the translated image.

### C. Local Discriminator

Masking the input of the discriminator may lead to the underrepresentation of some semantic classes. Therefore, we additionally train a local discriminator that operates on a batch of small patches sampled from the global crop. Our local discriminator is not masked but only sees patches where a certain amount of the semantic classes align. In our case, we sample patches with 1/8th the size of the global input crop where more than $50\%$ of the semantic classes align. We use our sampling procedure from Section III-A to sample these patches. Using small, partially aligned patches ensures that semantic classes are less underrepresented while maintaining content consistency.

### D. Feature-attentive Denormalization (FATE)

Spatially adaptive denormalization (SPADE) [56] fuses re-sized semantic segmentation maps as content into the generator stream. Feature-adaptive denormalization (FADE) [7] generalizes SPADE to features learned through a content stream. As shown in Figure 3, the normalized features $N(h)$ of the generator are modulated with the features $f$ of the content stream using the learned functions $\gamma$ and $\beta$ as follows:

$$\text{FADE}(h, f) = N(h) \circ \gamma(f) + \beta(f), \quad (4)$$

where $\gamma$ and $\beta$ are one-layer convolutions. This denormalization is applied in several layers of the generator. However, we argue that denormalization with content features is not always appropriate for transferring images to another domain because, as shown in [9], [57]–[59], image feature statistics contain not only content information but also style information. When
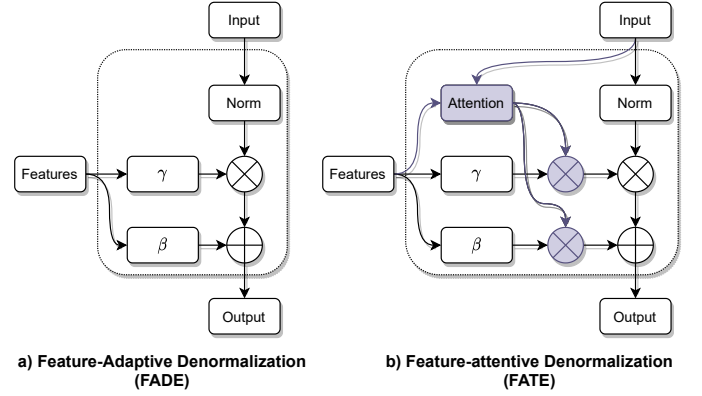


**a) Feature-Adaptive Denormalization (FADE)**   **b) Feature-attentive Denormalization (FATE)**

Fig. 3: **FADE and FATE.**

transferring feature statistics from the content stream of the source image to the generator stream, style information from the source image could affect the final image quality if the generator cannot ignore this information since the output image should mimic the style of the target domain. Therefore, we propose an additional attention mechanism to selectively incorporate statistics from the content stream into the generator stream. This allows the model to only fuse the statistical features from the source image into the generator stream that are useful for the target domain. As shown in Figure 3, this attention mechanism relies on the features of the content stream and the features of the generator stream and attends to the statistics $\gamma$ and $\beta$. With this attention mechanism, we can extend FADE to feature-attentive denormalization (FATE) as follows:

$$\text{FATE}(h, f) = N(h) \circ A(h, f) \circ \gamma(f) + A(h, f) \circ \beta(f), \quad (5)$$

where $A$ is the attention mechanism and $A(h, f)$ is the attention map for the statistics. We use a lightweight two-layer CNN with sigmoid activation in the last layer as the attention mechanism. More details can be found in Appendix A.

### E. Training Objective

Our training objective consists of three losses: a global masked adversarial loss $\mathcal{L}_{madv}^{global}$, a local adversarial loss $\mathcal{L}_{adv}^{local}$, and the perceptual loss $\mathcal{L}_{perc}$ used in [7]. We define the final training objective as follows:

$$\mathcal{L} = \lambda_{madv}^{global}\mathcal{L}_{madv}^{global} + \lambda_{adv}^{local}\mathcal{L}_{adv}^{local} + \lambda_{perc}\mathcal{L}_{perc}, \quad (6)$$

where we use a hinge loss to formulate the adversarial losses and $\lambda_{madv}^{global}$, $\lambda_{madv}^{local}$, $\lambda_{perc}$ are the corresponding loss weights.

## IV. EXPERIMENTS

### A. Experimental Settings

**Implementation details.** Our method is implemented in PyTorch 1.10.0 and trained on an A100 GPU (40 GB) with batch size 1. For training, we initialize all weights with the Xavier normal distribution [60] with a gain of 0.02 and use an Adam optimizer [61] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The initial learning rates of the generator and discriminators are set to 0.0001 and halved every $d_e$ epochs. Learning rate

decay is stopped after reaching a learning rate of 0.0000125. We formulate our adversarial objective with a hinge loss [62] and weight the individual parts of our loss function as follows: $\lambda_{madv}^{global} = 1.0, \lambda_{madv}^{local} = 1.0, \lambda_{perc} = 1.0$. In addition, we use a gradient penalty on target images [63], [64] with $\lambda_{rp} = 0.03$. The images of both domains are resized and cropped to the same size and randomly flipped before the sampling strategy is applied. In our experiments, we show that we achieve the best performance by cropping global patches of size 352×352. We crop local patches with 1/8th the size of the global crop (i.a., 44×44). The global discriminators are used on two scales. Crops are scaled down by a factor of two for the second scale. We train all our models for ∼ 400K iterations. Training a model takes 4-8 days, depending on the dataset, model, and crop size. We report all results as an average across five different runs. We refer to Appendix A for more details regarding the training and model. Our implementation is publicly available at https://github.com/BonifazStuhr/feamgan.

**Memory usage.** Our best model requires ∼25 GB of VRAM at training time and performs inference using ∼12 GB for an image of size 957×526. Our small model, with a slight performance decrease, runs on consumer graphic cards with ∼9 GB of VRAM at training time and performs inference using ∼8 GB for an image of size 957×526.

**Datasets.** We conduct experiments on four translation tasks across four datasets. For all datasets, we compute semantic segmentations with MSeg [55], which we use as a condition for our discriminator and to calculate the discriminator masks.

(1) *PFD* [65] consists of images of realistic virtual world gameplay. Each frame is annotated with pixel-wise semantic labels, which we use as additional input for our generator. We use the same subset as [2] to compare with recent work.

(2) *Viper* [66] consists of sequences of realistic virtual world gameplay. Each frame is annotated with different labels, where we use the pixel-wise semantic segmentations as additional input for our generator. Since Cityscapes does not contain night sequences, we remove them from the dataset.

(3) *Cityscapes* [67] consists of sequences of real street scenes from 50 different German cities. We use the sequences of the entire training set to train our models.

We use datasets (1-3) for the sim-to-real translation tasks *PFD→Cityscapes* and *Viper→Cityscapes*.

(4) *BDD100K* [68] is a large-scale driving dataset. We use subsets of the training and validation data for the following translation tasks: *Day→Night*, *Clear→Snowy*.

**Compared methods.** We compare our work with the following methods.

- Color Transfer (CT) [69] performs color correction by transferring statistical features in lαβ space from the target to the source image.
- MUNIT [15] achieves multimodal translation by recombining the content code of an image with a style code

sampled from the style space of the target domain. It is an extension of CycleGAN [17] and UNIT [22].
- CUT [31] uses a patchwise contrastive loss to achieve one-sided unsupervised image-to-image translation.
- TSIT [7] achieves one-sided translation by fusing features from the content stream into the generator on multiple scales using FADE and utilizing a perceptual loss between the translated and source images.
- QS-Attn [47] builds upon CUT [31] with an attention module that selects significant anchors for the contrastive loss instead of features from random locations of the image.
- EPE [2] relies on a variety of gbuffers as input. Techniques such as similarity cropping, utilizing segmentations for both domains generated by a robust segmentation model as input to the conditional discriminators, and small patch training are used to achieve content consistency.

Since EPE [2] provides inferred images of size 957×526 for the *PFD→Cityscapes* task, comparisons are performed on this resolution. For the *Viper→Cityscapes*, *Day→Night*, and *Clear→Snowy* tasks, we train the models using their official implementations. Furthermore, we retrain models as additional baselines for the *PFD→Cityscapes* task.

**Evaluation metrics.** Following prior work [2], we use the Fréchet Inception Distance (FID) [70], the Kernel Inception Distance (KID) [71], and the semantically aligned Kernel VGG Distance (sKVD) [2] to evaluate image translation quality quantitatively. The sKVD metric was introduced in [2] and improved over previous metrics for mismatched layouts in source and target data. In addition, we propose the class-specific Kernel VGG Distance (cKVD), where a robust segmentation model is used before the sKVD calculation to mask input crops by class (or category). Thereby, for each given class, all source and target image crops are filtered using their segmentations by erasing the pixels of all other classes. We select crops where more then 5% of the pixels belong to the respective class. Then, the sKVD is calculated class-wise on the filtered crops. Afterward, we can report the cKVD as an average over all classes or separately for each class to achieve a more fine-grained measurement. We follow [2] and use a crop size of 1/8 and sample source and target crop pairs with an similarity threshold of 0.5 between unmasked source and target segmentation crops. More information on the classes used in the cKVD metric can be found in Table IV of Appendix A. For the KID, sKVD, and cKVD metrics, we multiply the measurements by 1000 to improve the readability of results.

### B. Comparison to the State of the Art

We compare our models quantitatively and qualitatively with different baselines. First, we compare our results with EPE and the baselines provided by EPE [2]. Then, we train our own baselines on the four translation tasks for further comparison.

Fig. 4: **Qualitative comparison to EPE.** We compare our method with the provided inferred images of EPE [2].

TABLE I: **Quantitative comparison to the baselines provided by EPE.** We calculate all metrics on the provided inferred images of EPE and its baselines [2].

| Method | FID | KID | sKVD | cKVD | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | AVG | AVG$_{sp}$ | sky | ground | road | terrain | vegetation | building | roadside-obj. | person | vehicle | rest |
| ColorTransfer | 84.34 | 88.17 | 16.65 | 36.01 | 33.12 | 32.40 | **12.97** | 16.13 | 20.94 | 19.24 | 29.92 | 74.79 | 62.78 | 41.79 | **49.16** |
| MUNIT | 45.00 | 35.05 | 16.51 | 38.57 | 34.81 | 29.80 | 16.93 | 17.62 | 29.52 | 19.29 | **24.28** | 79.14 | 77.34 | 40.13 | 51.61 |
| CUT | 47.71 | 42.01 | 18.03 | 35.31 | 33.26 | **25.96** | 15.32 | 17.87 | **20.09** | 22.72 | 25.00 | 74.02 | **60.99** | 41.71 | 49.37 |
| EPE | 44.06 | 33.66 | 13.87 | **35.22** | **30.21** | 27.14 | 13.54 | **13.56** | 24.77 | 20.77 | 26.75 | **50.58** | 83.34 | 41.29 | 50.45 |
| FeaMGAN-S (ours) | **43.27** | **32.59** | **12.98** | 40.23 | 32.69 | 38.10 | 13.29 | 15.34 | 26.29 | 20.17 | 27.32 | 61.57 | 102.65 | 42.83 | 54.73 |
| FeaMGAN (ours) | **40.32** | **28.59** | **12.94** | 40.02 | 31.78 | 46.70 | 13.72 | 15.60 | 23.23 | **17.69** | 25.57 | 66.65 | 99.24 | **39.38** | 52.40 |

**Comparison to EPE.** A set of inferred images is provided for EPE and each of the baselines [2]. Therefore, we train our models on the same training set and use the inferred images from our best models for this comparison. We select our best models based on scores of various visual metrics and visual inspections of translated images. As shown in Figure 4 a) and b), our model relies solely on segmentation maps as additional input compared to EPE, which uses a variety of gbuffers. In addition, our model is trained with significantly fewer steps ($\sim$ 400K iterations) compared to EPE and the baselines (1M iterations). As shown in Table I, our model outperforms the baselines and EPE in all commonly used metrics (FID and KID) and the sKVD metric. More surprisingly, our small model, which can be trained on consumer GPUs, outperforms all baselines and EPE as well.

However, our cKVD metric shows that our models have difficulty with the person and sky classes. Therefore, the average cKVD values are high and become low when we remove both classes from the average calculation (AVG$_{sp}$). A possible reason for the weaker performance on the person

class is our masking procedure. Since the masking procedure requires overlapping samples in both domains, the person class is not seen frequently during training. This can lead to inconsistencies (a glow) around the person class, as seen in Figure 11 of our limitations. The masking procedure also leads to a drop in performance in the sky class, as seen in Table III of our ablation study. As shown in the first row of Figure 4 and the results of Figures 18 and 19 of Appendix A, our model translates larger structures, such as lane markings, more consistently, but fails to preserve some in-class characteristics from the source dataset. This is evident, for example, in the structure of translated streets and the corresponding cKVD value (road). As shown in the second row and Appendix A, EPE achieves visually superior modeling of the reflective properties of materials (e.g., the car) but suffers from inconsistencies (erased objects) regarding the vegetation, which can be seen in the palm trees and the corresponding cKVD value (vegetation). The superior modeling of reflective properties can be attributed to the availability of gbuffers (i.a., glossines) in EPE's input.

PFD→Cityscapes

Viper→Cityscapes

Day→Night

Clear→Snowy

| Input | MUNIT | CUT | TSIT | QS-Attn | FeaMGAN (ours) |

Fig. 5: **Qualitative comparison to prior work.** Models were trained using their official implementations. Randomly sampled results can be found in Figure 20 of Appendix A.

TABLE II: **Quantitative comparison to prior work.** Models were trained using their official implementations. Results are reported as the average across five runs. We refer to Table VII of Appendix A for an extended version of this table.

| Method | PFD→Cityscapes | | | | Viper→Cityscapes | | | | Day→Night | | | | Clear→Snowy | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FID | KID | sKVD | cKVD | FID | KID | sKVD | cKVD | FID | KID | sKVD | cKVD | FID | KID | sKVD | cKVD |
| Color Transfer | 91.01 | 94.82 | 18.16 | 50.87 | 89.30 | 83.51 | 20.20 | 51.23 | 125.90 | 140.60 | 32.58 | 56.52 | 46.85 | 19.44 | 14.91 | 42.89 |
| MUNIT | 40.36 | 29.98 | 14.99 | 43.24 | 47.96 | 30.35 | 14.14 | 59.62 | 42.53 | 31.83 | 15.02 | 50.83 | **44.74** | 17.48 | 11.65 | 48.10 |
| CUT | 49.55 | 44.25 | 16.85 | **37.53** | 60.35 | 49.48 | 16.80 | 51.02 | **34.36** | **20.54** | 10.16 | 53.55 | 46.03 | 15.70 | 14.71 | 43.91 |
| TSIT | **38.70** | **28.70** | **10.80** | 42.35 | **45.26** | **28.40** | **8.47** | 50.03 | 54.96 | 33.21 | 12.71 | 57.91 | 79.28 | 40.02 | 12.97 | 41.52 |
| QS-Attn | 49.41 | 42.87 | 14.01 | 38.57 | 55.62 | 39.31 | 12.99 | 63.22 | 46.67 | 21.47 | **7.58** | 52.02 | 60.91 | 18.85 | 14.19 | 44.00 |
| FeaMGAN-S (ours) | 45.16 | 34.93 | 13.87 | 40.50 | 52.79 | 35.92 | 14.34 | **45.38** | 70.40 | 51.30 | 14.68 | **46.66** | 57.93 | 16.24 | 11.88 | **38.28** |
| FeaMGAN (ours) | 46.12 | 36.56 | 13.69 | 41.19 | 51.56 | 34.63 | 14.01 | **47.21** | 66.39 | 46.96 | 13.14 | **46.88** | 56.78 | **14.77** | **11.36** | 41.72 |

By surpassing EPE in all commonly used quantitative metrics while maintaining content consistency, we are able to show that our model improves overall quantitative translation performance. However, our method has specific drawbacks that we discussed with the help of the cKVD metric and visual comparisons.

**Comparisons to retrained baselines.** We find that retraining the baselines with their original training setup for the PFD→Cityscapes task significantly improves their performance on commonly used metrics compared to the baselines provided by EPE, as can be seen in Table II. However, as shown in Figure 5 and the random results of Figure 20 of Appendix A, content-consistency problems remain. This indicates again that simply relying on commonly used metrics does not provide a complete picture if content consistency is taken into account. When qualitatively comparing our model to the baselines for the PFD→Cityscapes and Viper→Cityscapes tasks in Figure 5, we observe that our method significantly

reduces content inconsistencies. However, a limitation of our masking strategy are class boundary artifacts, which are particularly evident in the Day→Night translation task (Figure 11). Since masking allows our method to focus on specific classes, we achieve state-of-the-art performance for the Clear→Snowy translation task.

*C. Ablation Study*

**Effectiveness of masked discriminator.** As shown in Figure 6 and the random samples in Figure 22 of Appendix A, our masking strategy for the discriminator positively impacts content consistency. Without masking, inconsistencies occur that correlate with biases between the class distributions of the source and target domains. As shown in [2], the distributions of certain classes in the spatial image dimension vary greatly between the PFD dataset and the Cityscapes dataset. For example, trees in Cityscapes appear more frequently in the top half of the image, resulting in hallucinated trees when the images are translated without accounting for biases. In the

|     Input     |     Full     |   w/o Dis. Mask   |   w/o Local Dis.   |   w/ FADE w/o FATE   |

Fig. 6: **Qualitative ablations.** Results are selected from the best model. Randomly sampled results can be found in Figure 22 of Appendix A.

TABLE III: **Quantitative evaluation for ablation study.** Results are reported as the average across five runs. We refer to Table VIII of Appendix A for an extended version of this table.

| Method | FID | KID | sKVD | cKVD | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | AVG | sky | ground | road | terrain | vegetation | building | roadside-obj. | person | vehicle | rest |
| FeaMGan (Full) | 46.12 | 36.56 | 13.69 | 41.19 | 42.69 | 14.97 | 17.35 | 26.51 | **20.25** | 26.34 | 64.64 | 102.23 | **42.38** | 54.52 |
| w/o Dis. Mask | **37.10** | **25.88** | 14.73 | **39.65** | **26.70** | 15.81 | 16.65 | 31.02 | 22.97 | **25.39** | 67.01 | **93.78** | 44.23 | **52.91** |
| w/ FADE w/o FATE | 45.46 | 35.73 | **13.17** | 40.90 | 41.49 | 13.78 | 16.78 | 25.30 | 20.58 | 27.21 | **63.12** | 104.43 | 42.44 | 53.83 |
| w/ Random Crop | 47.88 | 38.48 | 13.37 | 40.18 | 39.88 | **12.90** | **14.65** | **25.09** | 21.89 | 27.32 | 64.32 | 98.81 | 43.08 | 53.86 |
| w/ VGG Crop | 51.23 | 42.46 | 13.56 | 40.62 | 40.32 | 13.38 | 15.67 | 26.47 | 21.09 | 27.28 | 65.23 | 99.61 | 43.19 | 53.94 |
| w/o Local Dis. | | | | | | | | | | | | | | |
| - w/ 256×256 Crop | 48.57 | 38.89 | **12.89** | 41.26 | 42.31 | 13.57 | 15.98 | **25.28** | 22.18 | **26.56** | **61.13** | 107.48 | 42.44 | 55.62 |
| - w/ 352×352 Crop | 47.26 | 37.75 | 14.38 | 39.30 | 34.44 | **13.09** | 15.84 | 25.83 | 21.50 | 27.20 | 61.24 | 98.25 | 42.24 | 53.38 |
| - w/ 464×464 Crop | **46.61** | **37.25** | 15.04 | **38.62** | **31.60** | 13.13 | **15.38** | 27.06 | 22.23 | 29.67 | 63.38 | **87.51** | 44.41 | 51.77 |
| - w/ 512×512 Crop | 55.89 | 49.12 | 15.94 | 39.35 | 36.48 | 14.68 | 16.06 | 26.87 | **19.61** | 27.37 | 62.40 | 98.90 | **40.32** | **50.86** |

first and second row of Figure 6, we show that our masking strategy (Full) prevents these inconsistencies in contrast to our model trained without masking (w/o Dis. Mask). However, as shown in Table III, this comes with a quantitative tradeoff in performance on commonly used metrics.

**Effectiveness of local discriminator.** We compare our model trained with a local discriminator (Full) to the model trained without a local discriminator (w/o Local Dis. 352x352). As shown in Figure 6, the local discriminator leads to an increase in quantitative performance. Furthermore, we show the qualitative effects of the local discriminator in Figure 6, where we observe a decrease in glowing objects and a significant decrease of erased objects in the translation. An example of a glowing object is the palm tree in row two of Figure 6. An example of erased objects are the missing houses in the background of the images from row three. In addition, small inconsistencies near object boundaries are reduced, as shown by the randomly sampled results in Figure 22 of Appendix A (e.g., the wheels of the car in row one and three). Overall, we can conclude that local discriminators can reduce local inconsistencies, which might arise from the

robust but not flawless segmentation maps used for masking.

**Effectiveness of segmentation-based sampling.** We compare our segmentation-based sampling method with random sampling and sampling based on VGG features. For the sampling strategy based on VGG features, we follow EPE [2] to calculate scores for 352×352 crops of the input images. Crops with a similarity score higher than 0.5 are selected for training. As shown in Table III, our segmentation-based sampling strategy (Full) slightly outperforms the other sampling strategies in overall translation performance.

**Effectiveness of FATE.** For each spatial point ("pixel") in the input feature map, our feature-attentive denormalization block selects the features in the feature dimension to be incorporated into the output stream of the generator by denormalization. We show the attention values of our feature-attentive denormalization block in Figure 9 by visualizing all attention values for a single feature across the entire feature map. Since a single feature represents a property of the input, a spatial pattern should emerge. This is expected especially in earlier layers, where the spatiality of the convolutional model's

| Input | 252×252 | 352×352 | 464×464 | 512×512 |

Fig. 7: **Qualitative ablation of crop sizes.** For each crop size, results are selected from the best model. Randomly sampled results can be found in Figure 21 of Appendix A.
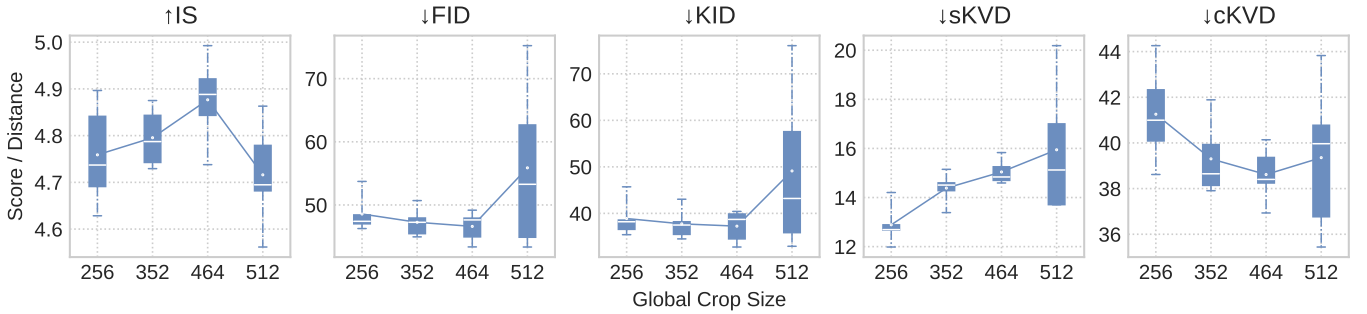


Fig. 8: **Quantitative ablation of crop sizes.**

feature map is best preserved. As shown in Figure 9, our attention mechanism learns to attend to features that correlate with a property. Examples are the shadows of a scene (row 1), cars and their lighting (row 2), and vegetation (row 3). In addition, we find increasingly more white feature maps in deeper layers. This can be interpreted positively as an indication that the learned content (source) features in deeper layers are important for the translation task and that more shallow content features of earlier layers are increasingly ignored. However, this can also be interpreted negatively and could indicate that our simple attention mechanism is not able to separate deeper features properly.

Comparing FATE to FADE, we find that FATE leads to a subtle increase in training instability, resulting in slightly worse average performance over the five runs per model. However, FATE also leads to our best models. Therefore, we select the FATE block as the standard configuration for our model. The deviation from the average values for all runs can be found in Table VIII of Appendix A. The slight increased instability suggests that the attention mechanism of FATE can be further improved.

**Effect of global crop size.** We successively increase the global crop size of the generator and discriminators from 256×256 to 512×512 and examine the effects on translation

performance. As shown in Figure 6, increasing the global crop size results in a better approximation of the target domain style. However, increasing the global crop size also leads to an increasing number of artifacts in the translated image. In Figure 8, we report the score of various metrics with respect to the global crop size. The commonly used metrics for measuring translation quality (IS, FID, and KID) show that translation quality increases steadily up to a global crop size of 464×464, after which the results become unstable. The cKVD metric also shows an increase in average performance up to a crop size of 464×464, mainly because translation quality for the underrepresented person class increases. This is intuitive since a larger crop size leads to a more frequent appearance of underrepresented classes during training. Furthermore, the sKVD metric shows a steady decline in consistency as the global crop size increases. Therefore, we choose a tradeoff between approximation of the target domain style, artifacts, and computational cost, and select 352×352 as the global crop size for our model.

## V. Conclusion

In this work, we have shown that content-based masking of the discriminator is sufficient to significantly reduce content inconsistencies that arise in unpaired image-to-image

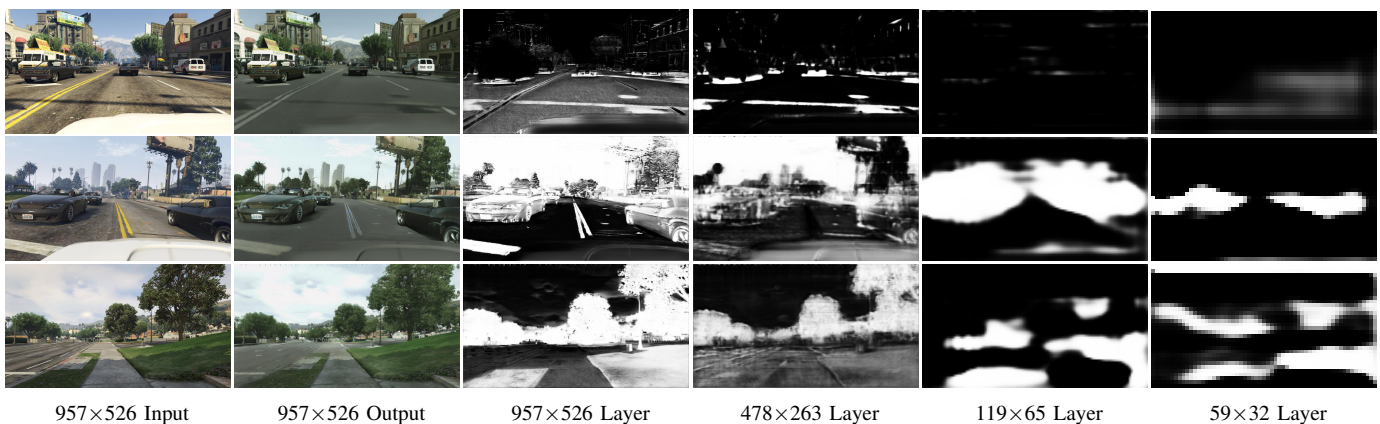| 957×526 Input | 957×526 Output | 957×526 Layer | 478×263 Layer | 119×65 Layer | 59×32 Layer |

Fig. 9: **FATE attention maps.** Results are selected from the best model.

translation. Furthermore, artifacts caused by the masking procedure can be significantly reduced by introducing a local discriminator that utilizes a segmentation-based similarity sampling technique. Moreover, our similarity sampling technique leads to a further increase in performance when applied to global input crops. We have also shown that our feature-based denormalization block is able to attend to specific content features, such as features of shadows, but can slightly increase training instability. In addition, we have proposed the cKVD metric to examine translation quality at the class or category level. In our experiments, we have found that these techniques lead to state-of-the-art performance on photo-realistic sim-to-real transfer and the translation of weather. Although our method performs well in Day→Night translation, the remaining limitations of our approach are especially evident in this task.

**Limitations.** We remark on limitations regarding the dataset, sampling, method, and implementation. Probably the most significant limitations are the complex public datasets currently available and in use, as they are not specifically designed for unpaired translation. Collection strategies and datasets that mitigate biases between source and target domains would be beneficial. Furthermore, our sampling strategy only works on an image basis and could be extended across the entire dataset to sample more significant pairs for training. Although our method works for large crops, there is still a crop size limit that must be taken into account when tuning the hyperparameters. In addition, our method for mitigating content inconsistencies depends on the segmentation model. In theory, the number of classes could be used to control how fine-grained the content consistency should be, which leads to flexibility but allows for errors depending on the segmentation quality. This can result in artifacts such as glowing objects, as shown in Figure 11. Intra-class inconsistencies that may arise from intra-class biases ignored by the loss, such as small textures, represent another problem. Intra-class inconsistencies are currently underexplored in unpaired image-to-image translation and are an interesting direction for future research. Finally, we would like to point out that the efficiency of our implementation

could be further improved. Apart from these limitations, our method achieves state-of-the-art performance in complex translation tasks while mitigating inconsistencies through a masking strategy that works by applying few tricks. Simple masking strategies have proven to be very successful in other fields. Therefore, we believe that masking strategies for unpaired image-to-image translation represent a promising direction for further research.

**Ethical and responsible use.** Considering the limitations of current methods, unpaired image-to-image translation methods should be trained and tested with care, especially for safety-critical domains like autonomous driving. A major concern is that it is often unclear or untested whether the transferred content can still be considered consistent for subsequent tasks in the target domain. Even though measures exist for content-consistent translation, they do not allow for the explainability of what exactly is being transferred and changed by the model on a fine-grained level. With our proposed cKVD metric we contribute to this field by allowing class-specific translation measurements - a direction that we hope is the right one. However, even if the content is categorically consistent at a high (class) level, subcategories (like parts of textures) may still be interchanged. At a lower level, content consistency and style consistency are intertwined (e.g., a yellow stop sign). Another privacy and security question is whether translation methods are (or will) be able to (indirectly) project sensitive information from the target domain to the translated images (e.g., exchange faces from simulation with faces of existing persons during the translation). A controllable (class-level and in-class-level) consistency method could help to resolve such issues.
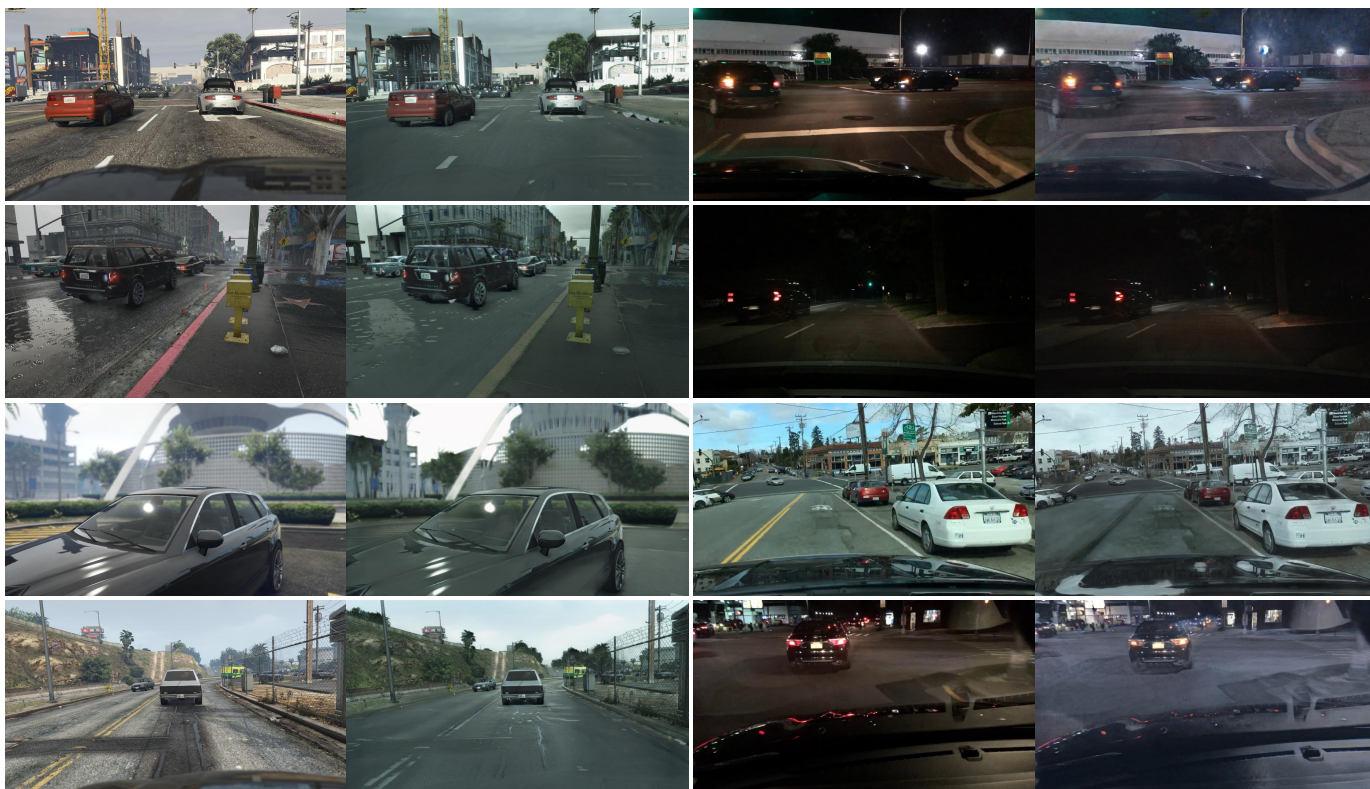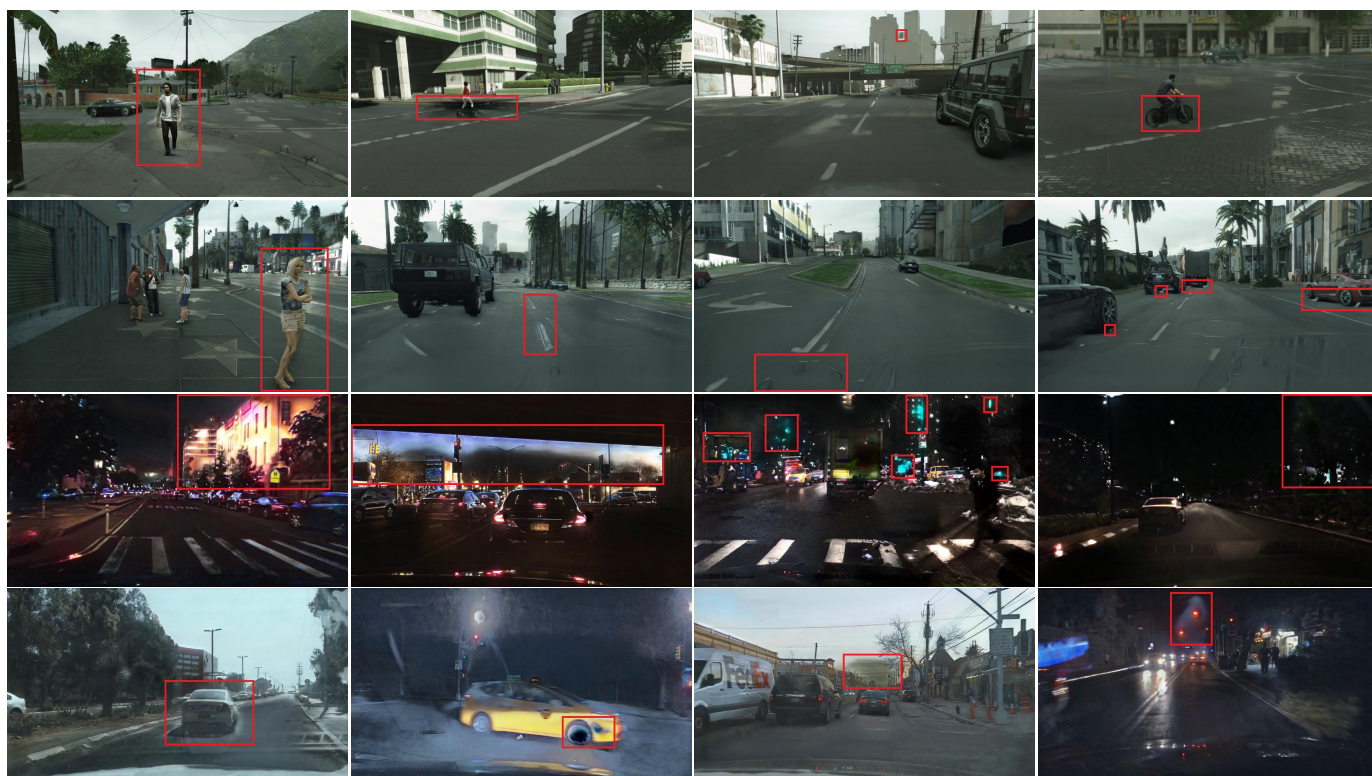
PFD→Cityscapes

Day→Night



Viper→Cityscapes

Clear→Snowy

Fig. 10: **Additional qualitative results.**

| Glowing objects | Intra-class inconsistencies | Minor hallucinations | Class boundary artifacts |

Fig. 11: **Limitations.**

REFERENCES

[1] F. Pizzati, P. Cerri, and R. de Charette, "Comogan: continuous model-guided image-to-image translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14 288–14 298.

[2] S. R. Richter, H. A. AlHaija, and V. Koltun, "Enhancing photorealism enhancement," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 2, pp. 1700–1715, 2022.

[3] Z. Jia, B. Yuan, K. Wang, H. Wu, D. Clifford, Z. Yuan, and H. Su, "Semantically robust unpaired image translation for data with unmatched semantics statistics," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14 273–14 283.

[4] Z. Hao, A. Mallya, S. Belongie, and M.-Y. Liu, "Gancraft: Unsupervised 3d neural rendering of minecraft worlds," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14 072–14 082.

[5] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, and T. Darrell, "Cycada: Cycle-consistent adversarial domain adaptation," in *International conference on machine learning*. Pmlr, 2018, pp. 1989–1998.

[6] S. Roy, A. Siarohin, E. Sangineto, N. Sebe, and E. Ricci, "Trigan: Image-to-image translation for multi-source domain adaptation," *Machine vision and applications*, vol. 32, pp. 1–12, 2021.

[7] L. Jiang, C. Zhang, M. Huang, C. Liu, J. Shi, and C. C. Loy, "Tsit: A simple and versatile framework for image-to-image translation," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*. Springer, 2020, pp. 206–222.

[8] S. Jeong, Y. Kim, E. Lee, and K. Sohn, "Memory-guided unsupervised image-to-image translation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 6558–6567.

[9] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1501–1510.

[10] Y. Yao, J. Ren, X. Xie, W. Liu, Y.-J. Liu, and J. Wang, "Attention-aware multi-stroke style transfer," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 1467–1475.

[11] J. Kim, M. Kim, H. Kang, and K. Lee, "U-gat-it: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation," *arXiv preprint arXiv:1907.10830*, 2019.

[12] C. Nederhood, N. Kolkin, D. Fu, and J. Salavon, "Harnessing the conditioning sensorium for improved image translation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6752–6761.

[13] J. Liang, H. Zeng, and L. Zhang, "High-resolution photorealistic image translation in real-time: A laplacian pyramid translation network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9392–9400.

[14] T. R. Shaham, M. Gharbi, R. Zhang, E. Shechtman, and T. Michaeli, "Spatially-adaptive pixelwise networks for fast image translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14 882–14 891.

[15] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz, "Multimodal unsupervised image-to-image translation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 172–189.

[16] C.-T. Lin, Y.-Y. Wu, P.-H. Hsu, and S.-H. Lai, "Multimodal structure-consistent image-to-image translation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 11 490–11 498.

[17] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.

[18] S. Benaim and L. Wolf, "One-sided unsupervised domain mapping," *Advances in neural information processing systems*, vol. 30, 2017.

[19] H. Fu, M. Gong, C. Wang, K. Batmanghelich, K. Zhang, and D. Tao, "Geometry-consistent generative adversarial networks for one-sided unsupervised domain mapping," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2427–2436.

[20] R. Zhang, T. Pfister, and J. Li, "Harmonic unpaired image-to-image translation," *arXiv preprint arXiv:1902.09727*, 2019.

[21] Y. Zhao, R. Wu, and H. Dong, "Unpaired image-to-image translation using adversarial consistency loss," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16*. Springer, 2020, pp. 800–815.

[22] M.-Y. Liu, T. Breuel, and J. Kautz, "Unsupervised image-to-image translation networks," *Advances in neural information processing systems*, vol. 30, 2017.

[23] O. Sendik, D. Cohen-Or, and D. Lischinski, "Crossnet: Latent cross-consistency for unpaired image translation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 3043–3051.

[24] Y. Yang, D. Lao, G. Sundaramoorthi, and S. Soatto, "Phase consistent ecological domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9011–9020.

[25] X. Liang, H. Zhang, L. Lin, and E. Xing, "Generative semantic manipulation with mask-contrasting gan," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 558–573.

[26] J. Theiss, J. Leverett, D. Kim, and A. Prakash, "Unpaired image translation via vector symbolic architectures," in *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXI*. Springer, 2022, pp. 17–32.

[27] C.-C. Kao, Y. Wang, J. Waltman, and P. Sen, "Patch-based image hallucination for super resolution with detail reconstruction from similar sample images," *IEEE Transactions on Multimedia*, vol. 22, no. 5, pp. 1139–1152, 2019.

[28] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014.

[29] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.

[30] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.

[31] T. Park, A. A. Efros, R. Zhang, and J.-Y. Zhu, "Contrastive learning for unpaired image-to-image translation," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16*. Springer, 2020, pp. 319–345.

[32] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*. Springer, 2016, pp. 694–711.

[33] X. Su, J. Song, C. Meng, and S. Ermon, "Dual diffusion implicit bridges for image-to-image translation," in *International Conference on Learning Representations*, 2022.

[34] M. Zhao, F. Bao, C. Li, and J. Zhu, "Egsde: Unpaired image-to-image translation via energy-guided stochastic differential equations," *arXiv preprint arXiv:2207.06635*, 2022.

[35] C. H. Wu and F. De la Torre, "Unifying diffusion models' latent space, with applications to cyclediffusion and guidance," *arXiv preprint arXiv:2210.05559*, 2022.

[36] Y. Taigman, A. Polyak, and L. Wolf, "Unsupervised cross-domain image generation," *arXiv preprint arXiv:1611.02200*, 2016.

[37] H. Wang, T. Shen, W. Zhang, L.-Y. Duan, and T. Mei, "Classes matter: A fine-grained adversarial approach to cross-domain semantic segmentation," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV*. Springer, 2020, pp. 642–659.

[38] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.

[39] X. Xie, J. Chen, Y. Li, L. Shen, K. Ma, and Y. Zheng, "Self-supervised cyclegan for object-preserving image-to-image domain adaptation," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*. Springer, 2020, pp. 498–513.

[40] L. Ma, X. Jia, S. Georgoulis, T. Tuytelaars, and L. Van Gool, "Exemplar guided unsupervised image-to-image translation with semantic consistency," *arXiv preprint arXiv:1805.11145*, 2018.

[41] X. Liu, G. Yin, J. Shao, X. Wang *et al.*, "Learning to predict layout-to-image conditional convolutions for semantic image synthesis," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[42] Y. Alami Mejjati, C. Richardt, J. Tompkin, D. Cosker, and K. I. Kim, "Unsupervised attention-guided image-to-image translation," *Advances in neural information processing systems*, vol. 31, 2018.

[43] H. Tang, H. Liu, D. Xu, P. H. Torr, and N. Sebe, "Attentiongan: Unpaired image-to-image translation using attention-guided generative adversarial networks," *IEEE transactions on neural networks and learning systems*, 2021.

[44] C. Yang, T. Kim, R. Wang, H. Peng, and C.-C. J. Kuo, "Show, attend, and translate: Unsupervised image translation with self-regularization and attention," *IEEE Transactions on Image Processing*, vol. 28, no. 10, pp. 4845–4856, 2019.

[45] L. Zhang, X. Chen, R. Dong, and K. Ma, "Region-aware knowledge distillation for efficient image-to-image translation," *arXiv preprint arXiv:2205.12451*, 2022.

[46] H. Tang, S. Bai, and N. Sebe, "Dual attention gans for semantic image synthesis," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 1994–2002.

[47] X. Hu, X. Zhou, Q. Huang, Z. Shi, L. Sun, and Q. Li, "Qs-attn: Query-selected attention for contrastive learning in i2i translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 291–18 300.

[48] H. Tang, D. Xu, N. Sebe, Y. Wang, J. J. Corso, and Y. Yan, "Multi-channel attention selection gan with cascaded semantic guidance for cross-view image translation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 2417–2426.

[49] G. Kwon and J. C. Ye, "Diagonal attention and style-based gan for content-style disentanglement in image generation and translation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 13 980–13 989.

[50] W. Liu, Z. Piao, Z. Tu, W. Luo, L. Ma, and S. Gao, "Liquid warping gan with attention: A unified framework for human image synthesis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 9, pp. 5114–5132, 2021.

[51] Y. Lin, Y. Wang, Y. Li, Y. Gao, Z. Wang, and L. Khan, "Attention-based spatial guidance for image-to-image translation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 816–825.

[52] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[53] D. Torbunov, Y. Huang, H. Yu, J. Huang, S. Yoo, M. Lin, B. Viren, and Y. Ren, "Uvcgan: Unet vision transformer cycle-consistent gan for unpaired image-to-image translation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 702–712.

[54] W. Zheng, Q. Li, G. Zhang, P. Wan, and Z. Wang, "Ittr: Unpaired image-to-image translation with transformers," *arXiv preprint arXiv:2203.16015*, 2022.

[55] J. Lambert, Z. Liu, O. Sener, J. Hays, and V. Koltun, "Mseg: A composite dataset for multi-domain semantic segmentation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2879–2888.

[56] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu, "Semantic image synthesis with spatially-adaptive normalization," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 2337–2346.

[57] L. A. Gatys, A. S. Ecker, and M. Bethge, "Image style transfer using convolutional neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2414–2423.

[58] C. Li and M. Wand, "Combining markov random fields and convolutional neural networks for image synthesis," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2479–2486.

[59] Y. Li, N. Wang, J. Liu, and X. Hou, "Demystifying neural style transfer," *arXiv preprint arXiv:1701.01036*, 2017.

[60] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010, pp. 249–256.

[61] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[62] J. H. Lim and J. C. Ye, "Geometric gan," *arXiv preprint arXiv:1705.02894*, 2017.

[63] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein gans," *Advances in neural information processing systems*, vol. 30, 2017.

[64] L. Mescheder, A. Geiger, and S. Nowozin, "Which training methods for gans do actually converge?" in *International conference on machine learning*. PMLR, 2018, pp. 3481–3490.

[65] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, "Playing for data: Ground truth from computer games," in *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*. Springer, 2016, pp. 102–118.

[66] S. R. Richter, Z. Hayder, and V. Koltun, "Playing for benchmarks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2213–2222.

[67] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3213–3223.

[68] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell, "Bdd100k: A diverse driving dataset for heterogeneous multitask learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2636–2645.

[69] E. Reinhard, M. Adhikhmin, B. Gooch, and P. Shirley, "Color transfer between images," *IEEE Computer graphics and applications*, vol. 21, no. 5, pp. 34–41, 2001.

[70] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," *Advances in neural information processing systems*, vol. 30, 2017.

[71] M. Bińkowski, D. J. Sutherland, M. Arbel, and A. Gretton, "Demystifying mmd gans," *arXiv preprint arXiv:1801.01401*, 2018.

## APPENDIX

### THE FEAMGAN ARCHITECTURE

**Generator**. As shown in Figure 12, our generator consists of a content stream encoder, a content stream, a generator stream encoder, and a generator stream. The content stream encoder shown in Figure 15 is utilized to create the initial features of the source image and condition. These initial features are the input to the content stream, which creates features for multiple levels with residual blocks. The statistics of these features are then integrated into the generator at multiple levels utilizing the residual FATE blocks shown in Figure 13. The generator stream utilizes the encoder shown in Figure 14 to create the initial latent from which the target image is generated. To further enforce content consistency, we do not use a variational autoencoder to obtain an deterministic latent. In addition, we found that utilizing additional residual blocks in the last layers of the generator stream improves performance, likely due to further refinement of the preceding upsampled features. We use spectral instance normalization for the residual blocks in the content stream and spectral batch normalization for the residual blocks in the generator stream. The convolutional layers in the generator stream encoder have the following numbers of filters: $[256, 512, 1024]$. The residual blocks in the generator have the following numbers of filters: $[1024, 1024, 1024, 512, 256, 128, 64, 64, 64, 64]$. The numbers of filters of the convolutional layers in the content streams encoder are $[64, 64]$. The numbers of filters in the content stream match those of the output of the preceding residual block in the generator stream at the respective level: $[64, 128, 256, 512, 1024, 1024, 1024, 1024]$. For all residual blocks, we use $3 \times 3$ convolutions and $1 \times 1$ convolutions for the skip connections. $\gamma$ and $\beta$ in the FATE and FADE blocks are created with $3 \times 3$ convolutions. Throughout the generator, we use a padding of 1 for the convolutions - we only downsample with strides and downsampling layers. We utilize the "nearest" upsampling and downsampling from Pytorch. For our small model, we halve the number of filters.

**Discriminator**. As shown in Figure 17, our discriminator consists of downsampling, upsampling, and prediction components. First, the input images of the source or target domain are downsampled via 5 stride 2 convolutions. We transform the output feature map of the last 4 downsampling convolutions with $1 \times 1$ convolutions. The last transformed feature map is used as input for the upsampling components, while the other transformed feature maps are added to the feature maps of the upsampling component for the receptive level. Then we utilize the feature maps of the 3 upsampling levels to create the final prediction on 3 levels. Thereby, we first apply a convolutional layer on the upsampled features. This convolution is followed by two convolutional layers: One is used to create the prediction feature map of depth 1, and the feature map of the other convolutional layer is multiplied by the segmentation map. The resulting segmentation feature map is then collapsed into depth 1 by adding the depth dimensions together. At last, the collapsed segmentation

TABLE IV: **cKVD class mapping.**

| cKVD Class | MSeg-Id(Name) |
|---|---|
| sky | 142(sky) |
| ground | 94(gravel), 95(platform), 97(railroad), 100(pavement-merged), 101(ground) |
| road | 98(road) |
| terrain | 102(terrain) |
| vegetation | 174(vegetation) |
| building | 31(tunnel), 32(bridge), 33(building-parent), 35(building), 36(ceiling-merged) |
| roadside-obj. | 130(streetlight), 131(road_barrier), 132(mailbox), 133(cctv_camera), 134(junction_box), 135(traffic_sign), 136(traffic_light), 137(fire_hydrant), 138(parking_meter), 139(bench), 140(bike_rack), 141(billboard) |
| person | 125(person), 126(rider_other), 127(bicyclist), 128(motorcyclist) |
| vehicle | 175(bicycle), 176(car), 177(autorickshaw), 178(motorcycle), 180(bus), 181(train), 182(truck), 183(trailer), 185(slow_wheeled_object) |
| rest | all other MSeg classes |

feature map is added to the prediction feature map to produce the final prediction. In this way, the discriminator is encouraged to produce class-specific predictions. We use spectral instance normalization for all convolutional layers. The $3 \times 3$ convolutions of the downsampling component have the following numbers of filters: $[64, 128, 256, 512, 512]$. The $1 \times 1$ convolutions of the downsampling component have the following numbers of filters: $[256, 256, 256, 256]$. The first convolutions in the prediction component have the following numbers of filters: $[128, 128, 128]$. The convolutions that are multiplied by the downsampled segmentation maps have the following numbers of filters: $[128, 128, 128]$. The convolution used to create the downsampled segmentation map has 128 filters. The convolutions to create the predictions have the following numbers of filters: $[1, 1, 1]$. Throughout the discriminator, we use a padding of 1 for the convolutions - we only downsample with strides and downsampling layers. We utilize the "bilinear" upsampling and downsampling from Pytorch. For our small model, we halve the number of filters.

### ADDITIONAL DATASET DETAILS

In Table V, we show additional details about the used datasets. Since we compare our method i.a. to the transferred PFD [65] images provided by EPE [2], we use a base image height of 526 throughout our experiments. The aspect ratio is preserved when resizing the input images. To match the input sizes of the images of both domains, we apply cropping to the image with the larger width if the image sizes of both domains do not align. The resulting images are randomly flipped before the sampling strategy is applied.

### ADDITIONAL TRAINING DETAILS

In Table VI, we show additional details about the hyperparameters used for training the four translation tasks. No tuning was performed for other translation tasks then PFD→Cityscapes besides adapting the learning rate schedule for the dataset lengths of these tasks.

**Content Stream**



**Generator Stream**

Fig. 12: **Generator architecture**. Arrows with dashed lines indicate connections at multiple levels between the two streams.



**FATE ResBlk**

Fig. 13: The FATE residual block used in the generator stream.



**Generator Stream Encoder**

Fig. 14: The generator stream encoder used to encode the input image and condition for the generator stream.



**Content Stream Encoder**

Fig. 15: The content stream encoder used to encode the input image and condition for the content stream.



**Attention**

Fig. 16: The attention module used in the FATE block to attend to the statistics of the features.

ADDITIONAL RESULTS

We show additional results of our experiments in Figures 18, 19, 20, 21, and 22. In Table VII, we report additional results from our cKVD metric and the stability of all results over five runs. Furthermore, we report the stability of all results from the ablation study in Table VIII. We note that the results for most baselines and for our method show non-negligible deviations in many tasks.

Fig. 17: **Discriminator architecture**. Arrows with dashed lines indicate connections at multiple levels between the two components.

TABLE V: **Additional details of the used datasets.**

| Dataset | Resolution | fps | Used Train/Val Data | Task | Input Resolution | Input Cropping |
|---------|-----------|-----|---------------------|------|------------------|----------------|
| PFD [65] | 1914×1052 | - | all images | *PFD→Cityscapes* | 957×526 | - |
| Viper [66] | 1920×1080 | ∼15 | all train/val data, but no night sequences | *Viper→Cityscapes* | 935×526 | - |
| Cityscapes [67] | 2048×1024 | 17 | all sequences of the train/val data | *PFD→Cityscapes* | 1.052×526 | 957×526 |
| | | | | *Viper→Cityscapes* | 1.052×526 | 935×526 |
| BDD100K [68] | 1280×720 | 30 | train: first 100k, val: first 40k | *Day→Night* | 935×526 | - |
| | | | train: first 50k, val: first 40k | *Clear→Snowy* | | |

TABLE VI: **Additional training details.**

| Task | Epochs | Schedule | Decay | Local Discriminator Batch Size |
|------|--------|----------|-------|--------------------------------|
| *PFD→Cityscapes* | 20 | half learning rate stepwise, learning rate ≥ 0.0000125 | after each 3rd epoch | 32 |
| *Viper→Cityscapes* | 5 | half learning rate stepwise, learning rate ≥ 0.0000125 | after each epoch | 32 |
| *Day→Night* | 5 | half learning rate stepwise, learning rate ≥ 0.0000125 | after each epoch | 32 |
| *Clear→Snowy* | 10 | half learning rate stepwise, learning rate ≥ 0.0000125 | after each epoch | 32 |

Input                                    EPE                                    FeaMGAN (ours)

Fig. 18: **Qualitative comparison to EPE.** We compare our method with the provided inferred images of EPE [2].

Input                               EPE                               FeaMGAN (ours)

Fig. 19: **Qualitative comparison to EPE.** We compare our method with the provided inferred images of EPE [2]. Results are randomly sampled from the best model.

PFD→Cityscapes

Viper→Cityscapes

Day→Night

Clear→Snowy

| Input | MUNIT | CUT | TSIT | QS-Attn | FeaMGAN (ours) |

Fig. 20: **Qualitative comparison to prior work.** Results are randomly sampled from the best model.

Input  252×252  352×352  464×464  512×512

Fig. 21: **Qualitative ablation of crop sizes.** For each crop size, results are randomly sampled from the best model.

Input  Full  w/o Dis. Mask  w/o Local Dis.  w/ FADE w/o FATE

Fig. 22: **Qualitative ablations.** Results are randomly sampled from the best model.

TABLE VII: **Extended quantitative comparison to prior work.** Models were trained using their official implementations. Results are reported as the average across five runs.

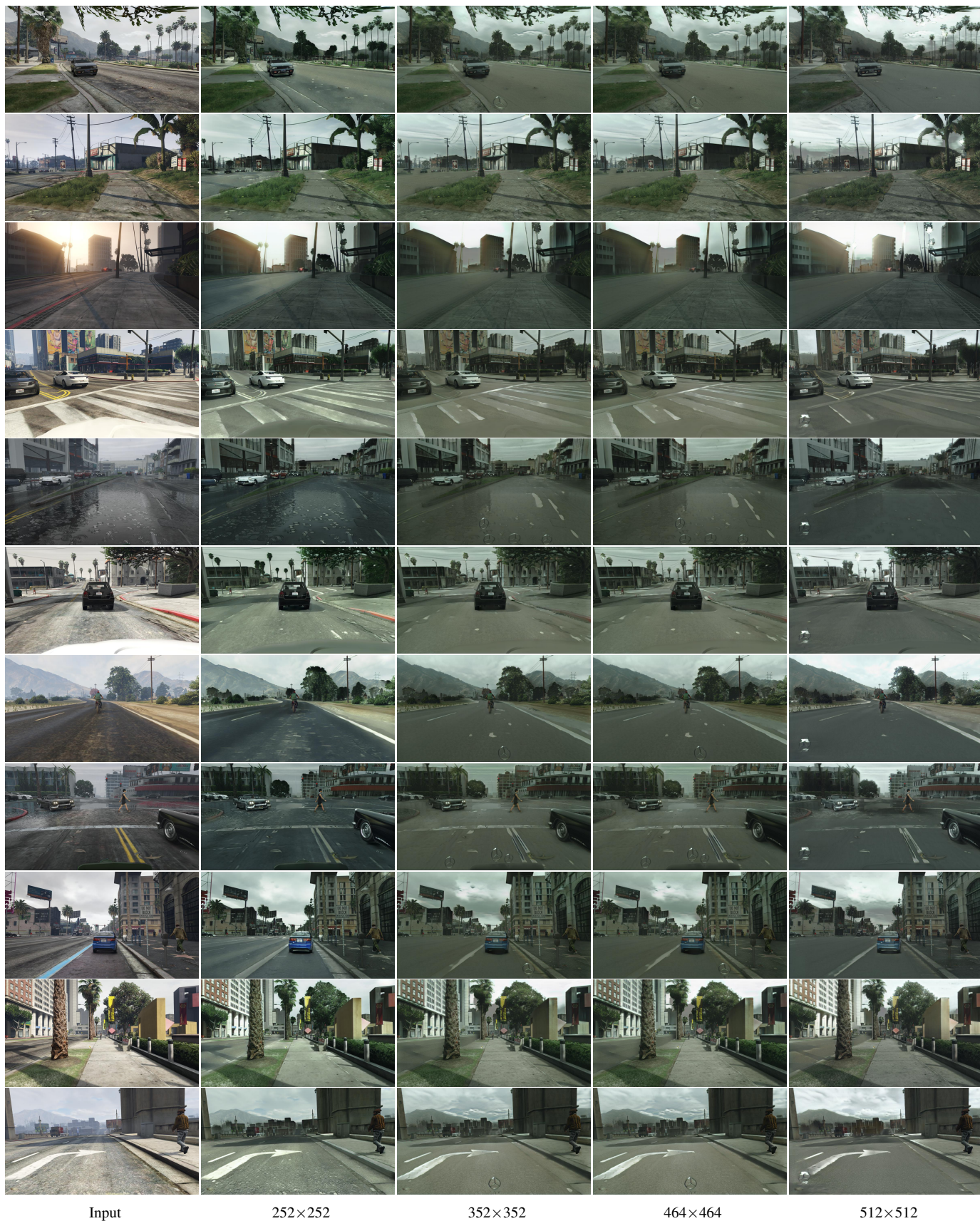| Method | FID | KID | sKVD | cKVD | | | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | AVG | sky | ground | road | terrain | vegetation | building | roadside-obj. | person | vehicle | rest |
| **PFD→Cityscapes** | | | | | | | | | | | | | | |
| Color Transfer | $91.01^{+0.05}_{-0.03}$ | $94.82^{+0.14}_{-0.13}$ | $18.16^{+0.20}_{-0.08}$ | $50.87^{+1.18}_{-0.36}$ | $58.05^{+2.28}_{-1.24}$ | $16.66^{+0.19}_{-0.26}$ | $16.38^{+0.14}_{-0.09}$ | $26.91^{+1.23}_{-1.29}$ | $28.18^{+0.44}_{-0.26}$ | $32.60^{+0.33}_{-0.45}$ | $58.36^{+5.56}_{-7.85}$ | $125.37^{+11.20}_{-7.66}$ | $55.12^{+1.38}_{-1.86}$ | $91.11^{+0.39}_{-0.84}$ |
| MUNIT | $40.36^{+1.22}_{-1.43}$ | $29.98^{+1.41}_{-1.08}$ | $14.99^{+0.18}_{-0.22}$ | $43.24^{+0.77}_{-0.71}$ | $37.92^{+0.24}_{-0.25}$ | $13.23^{+0.23}_{-0.23}$ | $14.33^{+0.09}_{-0.12}$ | $22.70^{+0.29}_{-0.75}$ | $24.97^{+0.35}_{-0.21}$ | $27.52^{+0.45}_{-0.30}$ | $58.24^{+3.87}_{-3.87}$ | $108.61^{+5.96}_{-7.85}$ | $45.35^{+0.82}_{-0.74}$ | $79.54^{+0.49}_{-0.85}$ |
| CUT | $49.55^{+3.63}_{-5.18}$ | $44.25^{+7.54}_{-4.99}$ | $16.85^{+1.28}_{-1.82}$ | $37.53^{+1.09}_{-1.13}$ | $28.76^{+0.98}_{-0.39}$ | $11.17^{+0.84}_{-0.27}$ | $13.92^{+0.40}_{-0.27}$ | $13.49^{+1.23}_{-1.21}$ | $24.20^{+1.39}_{-0.92}$ | $24.69^{+0.92}_{-0.86}$ | $57.45^{+4.33}_{-2.22}$ | $90.52^{+6.67}_{-0.16}$ | $40.92^{+2.77}_{-2.72}$ | $70.20^{+3.10}_{-0.99}$ |
| TSIT | $38.70^{+1.59}_{-1.16}$ | $28.70^{+1.81}_{-1.27}$ | $10.80^{+0.57}_{-0.29}$ | $42.35^{+0.89}_{-0.84}$ | $40.13^{+1.28}_{-1.14}$ | $13.74^{+0.78}_{-1.14}$ | $14.09^{+0.80}_{-0.78}$ | $23.48^{+1.17}_{-1.17}$ | $23.74^{+0.47}_{-0.36}$ | $25.76^{+0.16}_{-0.22}$ | $51.95^{+2.72}_{-1.88}$ | $107.98^{+5.61}_{-5.61}$ | $43.49^{+1.32}_{-1.01}$ | $79.13^{+2.72}_{-2.05}$ |
| QS-Attn | $49.42^{+5.71}_{-6.93}$ | $42.87^{+7.34}_{-10.03}$ | $14.01^{+0.34}_{-0.38}$ | $38.57^{+2.85}_{-1.68}$ | $29.50^{+2.40}_{-1.88}$ | $11.69^{+0.51}_{-0.69}$ | $13.92^{+0.68}_{-0.76}$ | $13.22^{+2.34}_{-1.23}$ | $23.99^{+1.32}_{-1.72}$ | $23.36^{+2.89}_{-1.56}$ | $57.88^{+2.89}_{-1.88}$ | $100.32^{+14.23}_{-6.18}$ | $40.77^{+3.32}_{-2.01}$ | $71.06^{+5.17}_{-3.29}$ |
| FeaMGan-S (ours) | $45.16^{+5.24}_{-3.23}$ | $34.93^{+6.92}_{-3.48}$ | $13.87^{+0.66}_{-0.89}$ | $40.50^{+2.83}_{-1.96}$ | $40.57^{+5.56}_{-4.83}$ | $13.32^{+2.35}_{-2.76}$ | $16.00^{+1.51}_{-2.76}$ | $24.55^{+2.88}_{-5.70}$ | $20.82^{+5.94}_{-4.83}$ | $27.54^{+1.07}_{-1.36}$ | $63.09^{+8.99}_{-4.94}$ | $102.53^{+1.58}_{-0.79}$ | $42.58^{+3.36}_{-3.01}$ | $53.99^{+5.41}_{-5.45}$ |
| FeaMGan (ours) | $46.12^{+4.60}_{-5.80}$ | $36.56^{+6.70}_{-7.96}$ | $13.69^{+1.13}_{-1.15}$ | $41.19^{+2.89}_{-2.81}$ | $42.69^{+4.00}_{-5.01}$ | $14.97^{+3.86}_{-3.23}$ | $17.35^{+5.09}_{-4.28}$ | $26.51^{+5.70}_{-3.27}$ | $20.25^{+2.88}_{-2.56}$ | $26.34^{+1.36}_{-0.77}$ | $64.64^{+4.94}_{-5.07}$ | $102.23^{+10.58}_{-10.91}$ | $42.38^{+2.91}_{-3.01}$ | $54.52^{+5.00}_{-3.09}$ |
| **Viper→Cityscapes** | | | | | | | | | | | | | | |
| Color Transfer | $89.30^{+0.06}_{-0.10}$ | $83.51^{+0.10}_{-0.68}$ | $20.20^{+0.32}_{-0.24}$ | $51.23^{+0.75}_{-1.87}$ | $65.74^{+1.33}_{-2.83}$ | $19.98^{+1.08}_{-0.87}$ | $16.87^{+0.16}_{-0.31}$ | $26.65^{+2.69}_{-2.23}$ | $28.79^{+0.45}_{-0.27}$ | $36.21^{+0.14}_{-0.24}$ | $41.97^{+1.52}_{-3.07}$ | $139.26^{+10.03}_{-15.46}$ | $57.10^{+1.04}_{-0.74}$ | $79.73^{+0.75}_{-0.82}$ |
| MUNIT | $47.96^{+0.52}_{-1.65}$ | $30.35^{+0.68}_{-1.29}$ | $14.14^{+0.09}_{-0.08}$ | $59.62^{+1.31}_{-2.34}$ | $46.44^{+2.59}_{-2.83}$ | $15.85^{+0.44}_{-0.71}$ | $14.11^{+0.31}_{-0.31}$ | $32.69^{+2.23}_{-3.43}$ | $25.75^{+0.30}_{-0.36}$ | $25.76^{+0.24}_{-0.24}$ | $39.99^{+3.65}_{-1.26}$ | $274.68^{+4.53}_{-23.29}$ | $46.64^{+1.75}_{-1.70}$ | $74.33^{+0.40}_{-0.52}$ |
| CUT | $60.35^{+6.50}_{-8.13}$ | $49.48^{+7.19}_{-10.15}$ | $16.80^{+1.11}_{-1.04}$ | $51.02^{+3.71}_{-4.32}$ | $34.79^{+6.87}_{-2.98}$ | $14.88^{+0.79}_{-0.78}$ | $16.80^{+2.50}_{-1.68}$ | $22.40^{+2.41}_{-2.33}$ | $22.91^{+1.81}_{-0.84}$ | $23.34^{+1.10}_{-1.04}$ | $45.00^{+5.23}_{-2.68}$ | $224.47^{+29.25}_{-28.29}$ | $42.29^{+2.56}_{-2.55}$ | $63.36^{+1.72}_{-3.17}$ |
| TSIT | $45.26^{+1.92}_{-1.39}$ | $28.40^{+2.55}_{-2.16}$ | $8.47^{+0.25}_{-0.26}$ | $50.03^{+3.06}_{-2.12}$ | $46.25^{+0.69}_{-0.93}$ | $14.46^{+2.11}_{-1.56}$ | $12.28^{+0.98}_{-0.97}$ | $31.95^{+4.96}_{-2.30}$ | $24.86^{+1.50}_{-1.43}$ | $24.91^{+1.26}_{-0.89}$ | $45.19^{+2.10}_{-2.35}$ | $184.05^{+18.06}_{-10.62}$ | $44.59^{+2.46}_{-1.55}$ | $71.72^{+3.90}_{-3.72}$ |
| QS-Attn | $55.62^{+12.05}_{-9.66}$ | $39.31^{+11.87}_{-2.55}$ | $12.99^{+1.27}_{-1.60}$ | $63.22^{+17.47}_{-13.74}$ | $36.44^{+15.38}_{-4.97}$ | $16.04^{+1.56}_{-1.26}$ | $15.25^{+1.10}_{-2.05}$ | $25.20^{+2.30}_{-2.30}$ | $26.09^{+1.43}_{-2.02}$ | $24.24^{+1.24}_{-0.89}$ | $46.54^{+1.63}_{-1.84}$ | $326.60^{+171.61}_{-128.90}$ | $46.44^{+3.78}_{-5.51}$ | $69.33^{+5.26}_{-5.06}$ |
| FeaMGan-S (ours) | $52.79^{+2.50}_{-2.79}$ | $35.92^{+3.88}_{-3.18}$ | $14.34^{+0.65}_{-0.73}$ | $45.38^{+1.53}_{-1.63}$ | $56.75^{+5.13}_{-8.76}$ | $18.51^{+1.49}_{-1.08}$ | $16.68^{+1.90}_{-3.38}$ | $42.85^{+1.59}_{-1.97}$ | $22.70^{+1.41}_{-1.40}$ | $26.82^{+0.74}_{-0.46}$ | $45.27^{+1.37}_{-1.11}$ | $130.76^{+6.27}_{-1.95}$ | $45.25^{+1.49}_{-3.29}$ | $48.19^{+2.45}_{-1.49}$ |
| FeaMGan (ours) | $51.56^{+2.17}_{-3.56}$ | $34.63^{+3.32}_{-5.48}$ | $14.01^{+0.58}_{-0.73}$ | $47.21^{+1.29}_{-1.10}$ | $58.87^{+3.48}_{-1.62}$ | $21.20^{+0.78}_{-0.93}$ | $18.03^{+1.38}_{-0.62}$ | $50.01^{+7.63}_{-3.72}$ | $23.55^{+3.08}_{-2.42}$ | $26.67^{+0.46}_{-0.66}$ | $45.55^{+1.11}_{-1.79}$ | $132.77^{+1.95}_{-3.46}$ | $45.13^{+3.72}_{-1.51}$ | $50.32^{+1.55}_{-1.24}$ |
| **Day→Night** | | | | | | | | | | | | | | |
| Color Transfer | $125.90^{+0.13}_{-0.10}$ | $140.60^{+0.10}_{-0.10}$ | $32.58^{+0.32}_{-0.52}$ | $56.52^{+1.76}_{-1.26}$ | $47.62^{+0.50}_{-0.78}$ | $27.41^{+1.37}_{-1.27}$ | $15.89^{+0.46}_{-0.23}$ | $32.60^{+1.76}_{-2.27}$ | $44.24^{+0.30}_{-0.25}$ | $32.61^{+0.68}_{-1.07}$ | $128.57^{+11.25}_{-13.18}$ | $108.52^{+8.17}_{-6.36}$ | $25.65^{+0.43}_{-0.37}$ | $102.06^{+0.39}_{-0.31}$ |
| MUNIT | $42.53^{+1.65}_{-1.27}$ | $31.83^{+1.73}_{-0.98}$ | $15.02^{+0.64}_{-0.65}$ | $50.83^{+1.25}_{-3.05}$ | $29.25^{+0.23}_{-0.30}$ | $28.00^{+0.52}_{-0.50}$ | $13.49^{+0.30}_{-0.64}$ | $36.57^{+0.53}_{-0.59}$ | $44.86^{+0.35}_{-0.69}$ | $24.96^{+0.69}_{-0.59}$ | $115.00^{+6.94}_{-5.17}$ | $101.70^{+4.06}_{-4.53}$ | $19.66^{+0.34}_{-0.33}$ | $94.82^{+1.31}_{-1.37}$ |
| CUT | $34.36^{+3.71}_{-6.12}$ | $20.54^{+4.81}_{-7.05}$ | $10.16^{+1.98}_{-1.14}$ | $53.55^{+3.05}_{-3.22}$ | $31.89^{+2.06}_{-2.32}$ | $27.44^{+1.58}_{-1.46}$ | $13.14^{+0.64}_{-0.73}$ | $40.93^{+6.70}_{-8.59}$ | $49.79^{+3.41}_{-2.78}$ | $25.52^{+1.69}_{-1.71}$ | $104.26^{+5.17}_{-10.28}$ | $122.50^{+11.90}_{-10.06}$ | $27.30^{+2.60}_{-2.87}$ | $92.76^{+0.52}_{-1.51}$ |
| TSIT | $54.979^{+6.12}_{-7.99}$ | $33.21^{+5.26}_{-6.21}$ | $12.71^{+5.77}_{-3.49}$ | $57.91^{+2.98}_{-2.28}$ | $36.27^{+1.39}_{-1.18}$ | $31.56^{+2.20}_{-2.21}$ | $16.93^{+1.07}_{-2.20}$ | $45.23^{+9.74}_{-4.86}$ | $54.82^{+4.89}_{-2.13}$ | $29.09^{+3.19}_{-1.65}$ | $143.47^{+14.01}_{-10.47}$ | $99.30^{+5.53}_{-3.53}$ | $27.43^{+2.98}_{-4.22}$ | $94.98^{+2.46}_{-2.60}$ |
| QS-Attn | $46.68^{+2.73}_{-2.03}$ | $21.47^{+2.34}_{-2.55}$ | $7.58^{+1.27}_{-1.77}$ | $52.02^{+4.14}_{-3.29}$ | $31.62^{+1.73}_{-1.66}$ | $26.73^{+2.64}_{-3.41}$ | $13.26^{+0.99}_{-0.92}$ | $38.25^{+5.42}_{-3.84}$ | $47.26^{+2.05}_{-3.13}$ | $25.42^{+2.05}_{-1.80}$ | $100.84^{+11.90}_{-7.85}$ | $123.79^{+18.23}_{-14.72}$ | $26.67^{+4.28}_{-3.72}$ | $86.39^{+3.62}_{-4.37}$ |
| FeaMGan-S (ours) | $70.40^{+15.29}_{-4.76}$ | $51.30^{+21.06}_{-6.07}$ | $14.68^{+3.45}_{-3.84}$ | $46.66^{+2.63}_{-2.56}$ | $30.35^{+1.05}_{-2.56}$ | $35.47^{+3.71}_{-1.75}$ | $17.26^{+1.46}_{-1.98}$ | $47.29^{+10.32}_{-8.77}$ | $27.12^{+1.14}_{-1.10}$ | $25.22^{+1.38}_{-2.72}$ | $116.25^{+4.47}_{-4.70}$ | $70.75^{+4.87}_{-5.91}$ | $19.29^{+0.52}_{-0.64}$ | $77.60^{+4.37}_{-2.86}$ |
| FeaMGan (ours) | $66.39^{+4.39}_{-8.43}$ | $46.96^{+6.07}_{-10.41}$ | $13.14^{+2.01}_{-2.01}$ | $46.88^{+1.53}_{-2.83}$ | $29.72^{+2.84}_{-1.42}$ | $35.94^{+1.17}_{-1.77}$ | $17.48^{+1.44}_{-1.44}$ | $49.78^{+7.45}_{-7.78}$ | $28.78^{+1.94}_{-3.03}$ | $25.65^{+0.71}_{-0.61}$ | $115.66^{+4.72}_{-5.03}$ | $70.94^{+9.41}_{-11.48}$ | $19.23^{+0.84}_{-0.78}$ | $75.57^{+4.73}_{-3.73}$ |
| **Clear→Snowy** | | | | | | | | | | | | | | |
| Color Transfer | $46.85^{+0.12}_{-0.18}$ | $19.44^{+0.43}_{-1.43}$ | $14.91^{+0.86}_{-2.97}$ | $42.89^{+13.81}_{-3.96}$ | $25.78^{+1.19}_{-2.03}$ | $22.99^{+3.22}_{-2.15}$ | $16.01^{+0.35}_{-0.21}$ | $21.54^{+1.72}_{-1.60}$ | $41.13^{+7.40}_{-0.58}$ | $24.20^{+1.99}_{-0.13}$ | $57.67^{+18.52}_{-11.79}$ | $128.26^{+94.11}_{-11.79}$ | $25.95^{+7.36}_{-2.36}$ | $65.39^{+5.38}_{-1.91}$ |
| MUNIT | $44.74^{+1.23}_{-0.79}$ | $17.48^{+0.59}_{-0.86}$ | $11.65^{+0.34}_{-0.22}$ | $48.10^{+0.49}_{-0.73}$ | $28.47^{+0.78}_{-0.78}$ | $25.64^{+0.30}_{-0.46}$ | $15.27^{+0.21}_{-0.48}$ | $26.21^{+0.60}_{-0.63}$ | $40.31^{+0.58}_{-0.47}$ | $24.26^{+0.13}_{-0.39}$ | $101.98^{+1.73}_{-4.77}$ | $116.14^{+3.89}_{-4.34}$ | $21.63^{+0.18}_{-0.41}$ | $81.08^{+0.31}_{-0.63}$ |
| CUT | $46.03^{+1.08}_{-0.85}$ | $15.70^{+0.77}_{-0.94}$ | $14.71^{+1.15}_{-0.50}$ | $43.91^{+0.97}_{-0.78}$ | $26.74^{+0.50}_{-0.63}$ | $21.96^{+0.96}_{-0.75}$ | $13.15^{+0.31}_{-0.50}$ | $21.49^{+1.02}_{-0.94}$ | $35.20^{+0.47}_{-0.63}$ | $25.31^{+0.59}_{-0.39}$ | $76.67^{+4.18}_{-6.21}$ | $119.13^{+16.19}_{-10.00}$ | $23.91^{+0.59}_{-1.18}$ | $75.51^{+1.17}_{-0.60}$ |
| TSIT | $79.29^{+5.08}_{-6.69}$ | $40.02^{+7.17}_{-7.14}$ | $12.97^{+0.37}_{-0.51}$ | $41.52^{+3.31}_{-2.59}$ | $28.02^{+2.81}_{-1.34}$ | $22.72^{+2.06}_{-2.52}$ | $14.32^{+0.64}_{-0.48}$ | $18.92^{+2.42}_{-2.13}$ | $34.54^{+2.08}_{-2.20}$ | $23.02^{+0.69}_{-0.66}$ | $72.13^{+7.37}_{-4.70}$ | $104.05^{+12.49}_{-12.64}$ | $21.64^{+2.06}_{-1.56}$ | $75.84^{+4.73}_{-5.53}$ |
| QS-Attn | $60.91^{+0.79}_{-1.02}$ | $18.85^{+1.05}_{-1.36}$ | $14.19^{+1.70}_{-1.01}$ | $44.00^{+1.95}_{-1.80}$ | $25.60^{+0.50}_{-0.15}$ | $22.04^{+1.20}_{-1.41}$ | $13.24^{+0.17}_{-0.71}$ | $22.71^{+0.46}_{-0.71}$ | $36.02^{+2.20}_{-1.42}$ | $26.45^{+1.75}_{-1.75}$ | $78.58^{+13.76}_{-4.53}$ | $114.07^{+13.76}_{-8.81}$ | $25.17^{+1.42}_{-2.29}$ | $76.08^{+1.84}_{-2.29}$ |
| FeaMGan-S (ours) | $57.93^{+1.37}_{-0.81}$ | $16.24^{+1.19}_{-0.98}$ | $11.88^{+0.55}_{-0.21}$ | $38.28^{+2.86}_{-2.61}$ | $22.69^{+0.67}_{-1.76}$ | $25.71^{+1.14}_{-3.01}$ | $15.82^{+0.82}_{-1.74}$ | $37.47^{+3.98}_{-4.13}$ | $25.94^{+1.50}_{-2.95}$ | $21.80^{+0.53}_{-0.61}$ | $75.32^{+4.53}_{-3.72}$ | $81.47^{+22.52}_{-14.46}$ | $19.10^{+0.55}_{-0.30}$ | $57.46^{+2.18}_{-3.99}$ |
| FeaMGan (ours) | $56.78^{+0.81}_{-0.32}$ | $14.77^{+0.52}_{-1.59}$ | $11.36^{+0.21}_{-0.43}$ | $41.72^{+2.61}_{-1.76}$ | $22.71^{+1.17}_{-1.09}$ | $26.64^{+3.01}_{-2.07}$ | $16.19^{+1.74}_{-1.13}$ | $38.00^{+3.46}_{-2.77}$ | $27.78^{+2.95}_{-2.21}$ | $21.41^{+0.61}_{-0.69}$ | $79.35^{+3.72}_{-1.97}$ | $105.08^{+35.80}_{-25.18}$ | $19.50^{+0.67}_{-1.04}$ | $60.59^{+3.99}_{-1.90}$ |

TABLE VIII: **Extended quantitative evaluation for ablation study.** Results are reported as the average across five runs.

| Method | FID | KID | sKVD | cKVD | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | AVG | sky | ground | road | terrain | vegetation | building | roadside-obj. | person | vehicle | rest |
| FeaMGan (Full) | $46.11^{+4.60}_{-5.80}$ | $36.66^{+6.70}_{-7.96}$ | $13.69^{+1.13}_{-1.15}$ | $41.19^{+2.89}_{-2.81}$ | $42.69^{+4.00}_{-5.01}$ | $14.97^{+3.86}_{-3.23}$ | $17.35^{+5.09}_{-4.28}$ | $26.51^{+5.70}_{-3.27}$ | $\mathbf{20.25}^{+2.88}_{-2.56}$ | $26.34^{+1.36}_{-0.77}$ | $64.64^{+4.94}_{-5.07}$ | $102.23^{+10.58}_{-10.91}$ | $\mathbf{42.38}^{+2.91}_{-3.01}$ | $54.52^{+5.00}_{-3.09}$ |
| w/o Dis. Mask | $\mathbf{37.10}^{+3.03}_{-3.60}$ | $\mathbf{25.88}^{+3.60}_{-8.65}$ | $14.73^{+1.27}_{-1.66}$ | $\mathbf{39.65}^{+4.73}_{-4.16}$ | $\mathbf{26.70}^{+4.09}_{-5.99}$ | $15.81^{+3.67}_{-3.22}$ | $16.65^{+4.43}_{-4.99}$ | $31.02^{+11.19}_{-8.23}$ | $22.97^{+4.15}_{-6.89}$ | $\mathbf{25.39}^{+2.75}_{-1.24}$ | $67.01^{+6.63}_{-7.99}$ | $\mathbf{93.78}^{+10.71}_{-7.90}$ | $44.23^{+2.88}_{-4.95}$ | $\mathbf{52.91}^{+5.34}_{-5.61}$ |
| w/ FADE w/o FATE | $45.46^{+2.68}_{-2.65}$ | $35.73^{+4.44}_{-3.53}$ | $\mathbf{13.17}^{+0.54}_{-0.60}$ | $40.90^{+4.40}_{-2.04}$ | $41.49^{+8.66}_{-8.86}$ | $13.78^{+3.22}_{-1.70}$ | $16.78^{+2.30}_{-1.62}$ | $25.30^{+2.41}_{-2.12}$ | $20.58^{+1.88}_{-5.25}$ | $27.21^{+0.68}_{-1.93}$ | $\mathbf{63.12}^{+3.77}_{-5.25}$ | $104.43^{+13.22}_{-7.73}$ | $42.44^{+0.91}_{-1.27}$ | $53.83^{+1.96}_{-5.46}$ |
| w/ Random Crop | $47.88^{+5.10}_{-3.82}$ | $38.48^{+7.27}_{-4.14}$ | $13.37^{+0.61}_{-1.14}$ | $40.18^{+1.55}_{-1.66}$ | $39.88^{+3.92}_{-7.65}$ | $\mathbf{12.90}^{+0.94}_{-0.68}$ | $\mathbf{14.65}^{+1.62}_{-2.00}$ | $\mathbf{25.09}^{+2.01}_{-2.93}$ | $21.89^{+5.21}_{-2.05}$ | $27.32^{+2.05}_{-1.87}$ | $64.32^{+2.86}_{-1.79}$ | $98.81^{+7.63}_{-5.61}$ | $43.08^{+3.81}_{-2.90}$ | $53.86^{+1.66}_{-1.56}$ |
| w/ VGG Crop | $51.23^{+5.12}_{-2.11}$ | $42.46^{+6.65}_{-3.42}$ | $13.56^{+0.78}_{-0.81}$ | $40.62^{+2.22}_{-1.99}$ | $40.32^{+4.36}_{-5.72}$ | $13.3^{+2.06}_{-1.83}$ | $15.67^{+2.15}_{-2.14}$ | $26.47^{+4.07}_{-1.25}$ | $21.09^{+1.61}_{-2.41}$ | $27.28^{+0.98}_{-1.46}$ | $65.23^{+8.64}_{-3.75}$ | $99.61^{+9.27}_{-7.66}$ | $43.19^{+2.16}_{-3.12}$ | $53.94^{+3.81}_{-5.04}$ |
| w/o Local Dis. | | | | | | | | | | | | | | |
| - w/ 256×256 Crop | $48.57^{+5.16}_{-2.30}$ | $38.89^{+6.82}_{-3.47}$ | $\mathbf{12.89}^{+1.32}_{-0.90}$ | $41.26^{+3.00}_{-2.64}$ | $42.31^{+2.92}_{-3.88}$ | $13.57^{+1.65}_{-2.15}$ | $15.98^{+1.44}_{-1.58}$ | $\mathbf{25.28}^{+4.26}_{-2.75}$ | $22.18^{+8.57}_{-4.62}$ | $\mathbf{26.56}^{+2.06}_{-1.46}$ | $\mathbf{61.13}^{+2.92}_{-2.21}$ | $107.48^{+6.40}_{-11.62}$ | $42.44^{+4.20}_{-2.30}$ | $55.62^{+5.79}_{-4.61}$ |
| - w/ 352×352 Crop | $47.26^{+3.44}_{-2.31}$ | $37.75^{+5.30}_{-3.23}$ | $14.38^{+0.76}_{-1.00}$ | $39.30^{+2.59}_{-1.40}$ | $34.44^{+7.70}_{-5.20}$ | $\mathbf{13.09}^{+1.62}_{-1.17}$ | $15.84^{+1.79}_{-1.04}$ | $25.83^{+1.56}_{-2.01}$ | $21.50^{+3.84}_{-1.63}$ | $27.20^{+1.24}_{-1.47}$ | $61.24^{+5.82}_{-1.24}$ | $98.25^{+4.82}_{-3.34}$ | $42.24^{+2.34}_{-1.79}$ | $53.38^{+3.71}_{-2.18}$ |
| - w/ 464×464 Crop | $\mathbf{46.61}^{+2.57}_{-3.25}$ | $\mathbf{37.25}^{+3.17}_{-4.48}$ | $15.04^{+0.79}_{-0.45}$ | $\mathbf{38.62}^{+1.52}_{-1.68}$ | $\mathbf{31.60}^{+4.89}_{-3.80}$ | $13.13^{+1.08}_{-0.99}$ | $\mathbf{15.38}^{+1.77}_{-1.93}$ | $27.06^{+1.37}_{-2.76}$ | $22.23^{+3.99}_{-3.40}$ | $29.67^{+0.67}_{-1.10}$ | $63.38^{+3.73}_{-3.93}$ | $\mathbf{87.51}^{+3.83}_{-4.99}$ | $44.41^{+2.76}_{-2.76}$ | $51.77^{+2.19}_{-3.06}$ |
| - w/ 512×512 Crop | $55.89^{+19.35}_{-12.54}$ | $49.12^{+26.93}_{-16.19}$ | $15.94^{+4.24}_{-2.24}$ | $39.35^{+4.47}_{-3.91}$ | $36.48^{+6.33}_{-12.08}$ | $14.68^{+2.82}_{-2.30}$ | $16.06^{+3.25}_{-3.12}$ | $26.87^{+8.35}_{-4.31}$ | $\mathbf{19.61}^{+4.76}_{-4.21}$ | $27.37^{+3.00}_{-2.34}$ | $62.40^{+3.93}_{-4.27}$ | $98.90^{+8.87}_{-5.05}$ | $\mathbf{40.32}^{+3.12}_{-4.21}$ | $\mathbf{50.86}^{+5.50}_{-7.19}$ |