

# Thermal Infrared Image Colorization for Nighttime Driving Scenes with Top-Down Guided Attention

Fuya Luo, Yunhan Li, Guang Zeng, Peng Peng, Gang Wang, and Yongjie Li, *Senior Member, IEEE*

**Abstract**—Benefitting from insensitivity to light and high penetration of foggy environments, infrared cameras are widely used for sensing in nighttime traffic scenes. However, the low contrast and lack of chromaticity of thermal infrared (TIR) images hinder the human interpretation and portability of high-level computer vision algorithms. Colorization to translate a nighttime TIR image into a daytime color (NTIR2DC) image may be a promising way to facilitate nighttime scene perception. Despite recent impressive advances in image translation, semantic encoding entanglement and geometric distortion in the NTIR2DC task remain under-addressed. Hence, we propose a top-down attention and gradient alignment based GAN, referred to as PearlGAN. A top-down guided attention module and an elaborate attentional loss are first designed to reduce the semantic encoding ambiguity during translation. Then, a structured gradient alignment loss is introduced to encourage edge consistency between the translated and input images. In addition, pixel-level annotation is carried out on a subset of FLIR and KAIST datasets to evaluate the semantic preservation performance of multiple translation methods. Furthermore, a new metric is devised to evaluate the geometric consistency in the translation process. Extensive experiments demonstrate the superiority of the proposed PearlGAN over other image translation methods for the NTIR2DC task. The source code and labeled segmentation masks will be available at <https://github.com/FuyaLuo/PearlGAN/>.

**Index Terms**—Thermal infrared image colorization, image-to-image translation, generative adversarial networks, guided attention, nighttime driving scenes perception.

## I. INTRODUCTION

**R**OBUST and reliable all-weather scene perception is essential for intelligent driving assistance systems. Nevertheless, imaging devices based on the visible spectrum are extremely sensitive to lighting conditions and fail in total darkness, which limits their practicality in severe weather and other environments with poor visibility. Thermal long-wave infrared (LWIR) cameras can capture infrared radiation (7.5–14  $\mu\text{m}$ ) emitted by objects with a temperature above absolute zero, which, unlike visible cameras, allows them to image low-light environments without the aid of an illumination source. However, thermal infrared (TIR) images captured by LWIR cameras usually have low contrast and ambiguous object boundaries. Furthermore, TIR images are in shortage of chrominance, which might hinder human interpretation [1], [2]

Fuya Luo, Yunhan Li, Guang Zeng, Peng Peng and Yongjie Li are with the MOE Key Laboratory for Neuroinformation, the School of Life Science and Technology, University of Electronic Science and Technology of China, Chengdu 610054, China. E-mail: luofuya1993@gmail.com, liyunhan@std.uestc.edu.cn, guangzeng0109@gmail.com, peng-panda.uestc@gmail.com, liyj@uestc.edu.cn. (*Corresponding authors: Yongjie Li and Gang Wang.*)

Gang Wang is with Center of Brain Sciences, Beijing Institute of Basic Medical Sciences, Beijing 100085. E-mail: g\_wang@foxmail.com.

and subsequent transplantation of high-level computer vision algorithms. Therefore, it is of great significance to transform a nighttime TIR (NTIR) image into a daytime color (DC) image, which not only helps the driver to quickly understand the surrounding environment and sense abnormalities when driving at night, but also reduces the annotation burden on NTIR image based computer vision tasks by utilizing massively annotated DC image datasets. In this paper, we aim to address the problem of NTIR image colorization, which is also called translation from NTIR image to DC image (abbreviated as NTIR2DC).

The aim of NTIR2DC is to map a single-channel TIR image to a 3-channel RGB image and brighten the scene without changing the semantics. In general, existing TIR image colorization approaches can be categorized into supervised [3], [4] and unsupervised [5] approaches. Due to the rapid changes in traffic scenarios, pixel-level registered infrared-visible image pairs are difficult to acquire, which limits the performance of the supervised approaches. In contrast, unsupervised image-to-image (I2I) translation offers a potential solution to this problem by enforcing the input data distributions of two domains to be similar without paired cross-domain images. Recently, the compelling success of Generative Adversarial Networks (GAN) [6] has brought new blood into unsupervised I2I translation. To loosen the requirement of pairwise training images, CycleGAN [7] introduced a cycle consistency loss, which attempts to preserve the original image after a cycle of translation and reverse translation. Nyberg *et al.* [5] utilized the CycleGAN model to realize unpaired infrared-visible image translation. MUNIT [8] was proposed to improve the diversity of synthetic images. Anoosheh *et al.* [9] utilized a night-to-day image translation model called ToDayGAN to improve image retrieval performance for localization tasks.

Although impressive results have been obtained through unsupervised I2I translation methods, the problem of geometrical distortion during translation is still under-addressed, and how to utilize contextual information to reduce local semantic ambiguity for NTIR image encoding is still under-explored. As shown in Fig. 1, both CycleGAN and MUNIT fail at generating plausible pedestrians and keeping the edges consistent with the original image. To address the above mentioned problems, we propose a new GAN model incorporated with top-down attention and gradient alignment based on ToDayGAN, referred to as PearlGAN.

Visual attention allows us to interact with our environment by selectively attending to the information that is relevant to our behavior. Studies [10] have shown that the spatial information about a scene extracted by the posterior parietal



Fig. 1. Visual comparison of detection results and geometric consistency. In the second column, the green dashed box is the zoomed-in result of the fusion between the corresponding region and the Canny edge of the original NTIR image, and the red dashed box is the zoomed-in result of the detection. Please zoom in to check more details on the content and quality.

cortex (PPC) forms the basis for feedback signals to highlight neural responses as early along the visual pathway as the primary visual cortex. Such feedback modulated filtering helps reduce information overload and enables effective continuous visual search by directing our attention to specific locations in the visual field [11]. Inspired by this attentional feedback, we devise a top-down guided attention (TDGA) module with an elaborate attentional loss to reduce local semantic ambiguity in infrared image coding. Specifically, the TDGA module first uses average pooling at different scales to extract spatial contextual information of different receptive fields. Then, the maximum receptive field extracted features serve as coarse global information (analogous to the scene information extracted by PPC), which directs the attention of different receptive field units to specific spatial locations from coarse to fine, so as to achieve hierarchical scene coding. The introduced attentional loss function includes attentional diversity loss and attentional cross-domain conditional similarity loss, which are used to encourage a diverse distribution of attentional maps at the spatial level and feature level, respectively. For the edge distortion problem, we design a structured gradient alignment loss to penalize edge shift or disappearance during the image translation process.

Given that image colorization is a multi-solution problem, a natural question is how to evaluate the accuracy of colorization without ground truth. Indeed, an acceptable colorized TIR image requires not only realistic features but also a rigorous preservation of scene content, which is difficult to fully assess using the FID score [12]. Therefore, we use a semantic segmentation model pre-trained on real DC images to evaluate the semantic preservation performance of various colorization methods. In addition, we propose a new metric called APCE (i.e., Average Precision of Canny Edges under multi-threshold conditions) to evaluate geometric consistency between the original and translated images.

To measure the semantic segmentation performance of the translated images, we perform pixel-level category annotation on a subset of NTIR images from the FLIR [13] and KAIST [14] datasets. Exhaustive experiments on these two

datasets show that the proposed method not only achieves plausible colorization, but also outperforms other state-of-the-art translation methods in terms of image content and geometry preservation. The main contributions of this study are summarized as follows:

- We design a top-down guided attention module and a corresponding attentional loss to achieve hierarchical attention distribution and reduce local semantic ambiguity of image encoding using contextual information.
- We introduce a structured gradient alignment loss to reduce edge distortion in the NTIR2DC task.
- We annotate a subset of FLIR and KAIST datasets with pixel-level category labels, which may catalyze research on the colorization and semantic segmentation of NTIR images.
- To the best of our knowledge, we are the first to propose evaluation metrics to assess the semantic and edge preservation of NTIR2DC methods.
- Extensive experiments on the NTIR2DC task show that the proposed model significantly outperforms other image translation methods in terms of semantic preservation and edge consistency.

In the remainder of this paper, Section II summarizes related work about TIR image colorization and I2I translation. Section III introduces the architecture of the proposed PearlGAN. Section IV presents our experiments on FLIR and KAIST datasets. Section V draws the conclusions of our work.

## II. RELATED WORK

In this section, we briefly review previous work on TIR image colorization, unpaired image-to-image translation and feature aggregation with attention pyramid.

### A. TIR Image Colorization

Recently, witnessing the success of deep learning based approaches in various computer vision tasks [15], [16], [17], an increasing number of researchers have focused their efforts on the area of TIR image colorization, which can be categorized into supervised [4], [3], [18], [19], [20], [21], [22], [23] and unsupervised [5] methods. The supervised colorization methods require paired infrared and visible images, and they colorize TIR images by minimizing the distance between the synthesized image and the corresponding RGB image. For example, Berg *et al.* [4] used separate luminance and chrominance loss to constrain the mapping of infrared to visible RGB images. [21] and [22] combined pixel-level content loss and adversarial loss to realize the colorization of thermal infrared images. In order to enhance edge information in the image translation process, Wang *et al.* [18], [23] jointly encoded infrared images and their Canny edge maps.

However, due to the rapid changes in specific scenarios (e.g., traffic scenes), it is extremely difficult to collect TIR and visible image pairs with perfect pixel to pixel correspondence, which limits the practicality of supervised methods. Without requiring paired samples, unsupervised colorization methods usually use GAN models to minimize the distance between the distribution of the translated image and the real RGB image.

For example, Nyberg *et al.* [5] utilized the CycleGAN [7] model to realize unpaired infrared-visible image translation. In general, despite the impressive results obtained with the previous methods, semantic encoding entanglement and geometric distortion in the NTIR2DC task is still under-addressed.

### B. Unpaired Image-to-Image Translation

Unpaired image-to-image translation aims to use unpaired training samples to learn the mapping between two different image domains. Zhu *et al.* [7] made the earliest effort to get rid of aligned image pairs by using cycle consistency loss, leading to the recent surge of interest in methods for unpaired image-to-image translation [24], [25], [26]. Based on the shared-latent space assumption, Liu *et al.* [24] presented a general framework named UNIT to further extend the cycle-consistency constraint. Subsequently, MUNIT [8] and DRIT++ [27] were proposed to improve the diversity of synthesized images by learning a disentangled representation with a domain-invariant content space and a domain-specific style space. Recent works further boosted the generation performance of night-to-day image translation by introducing multiple discriminators [9] or decoders [28]. Compared with the implicit edge consistency constraint using cycle consistency loss, the structured gradient alignment loss proposed in this work explicitly constrains the gradient structure of the synthesized image, to avoid the possible edge misalignment repair in the inverse mapping process.

### C. Feature Aggregation with Attention Pyramid

Recently, due to its superiority in multi-scale feature selection, a strategy that merges features with an attention pyramid has been widely used in salient object detection [29], [30], [31], face recognition [32] and video classification [33]. For example, Wang *et al.* [29] proposed a pyramid attention module, which extends the regular attention mechanisms with multiscale information to improve saliency representation. Ni *et al.* [34] proposed a Double Attention Module to extract multi-scale attentive features, which were fused through a Pyramid Upsampling Module. To extract the most discriminative semantic regions for domain adaptation, Li *et al.* [35] adopted a task-oriented guided spatial attention pyramid learning strategy to aggregate hierarchical semantic information in feature maps of the pyramid. Hu *et al.* [36] constructed a pyramid of local self-attention blocks to achieve localization of image manipulations.

The module most similar to the proposed TDGA module might be the Cascaded Pyramid Attention (CPA) module [31], as both progressively estimate high-resolution attention maps from coarse to fine. However, different from the CPA module that uses only the largest scale attention map to select input features, the TDGA module uses attention maps of different scales to selectively weight the corresponding groups of feature maps. In terms of purpose, the TDGA module focuses on the hierarchical attention of the whole scene and the disentanglement of spatial features, whereas the CPA module aims to strengthen the regional features of specific objects.

## III. PROPOSED METHOD

In this section, we first present the overview of the proposed PearlGAN. Subsequently, we briefly explain the ToDayGAN model [9] as our baseline and its variants for TIR image colorization. Then the details of the proposed TDGA module are described. Next, we explicate the elaborate attentional loss, including attentional diversity (AD) loss and attentional cross-domain conditional similarity (ACCS) loss. After that, the structured gradient alignment (SGA) loss to enable edge consistency during translation is explained. Finally, we illustrate the total loss of the proposed PearlGAN.

### A. Overview and Problem Formulation

The overall framework is shown in Fig. 2. We choose the ToDayGAN [9] model as the baseline model, which consists of a pair of generators and discriminators with a total objective function including adversarial loss  $\mathcal{L}_{adv}$  and cycle-consistency loss  $\mathcal{L}_{cyc}$ . The generator consists of an encoder and a decoder. Due to its limitations on the NTIR2DC task, we first adapt the ToDayGAN model using the existing loss functions and modules to obtain a new model called ToDayGAN-TIR. Based on the ToDayGAN-TIR model, we introduce a novel TDGA module and an attentional loss to reduce the coding ambiguity of NTIR images, and a SGA loss  $\mathcal{L}_{sga}$  to reduce the edge distortion of translation results. The final model obtained is called PearlGAN. Attentional loss consists of an AD loss  $\mathcal{L}_{att}^{div}$  and an ACCS loss  $\mathcal{L}_{att}^{ccs}$ .

In the rest of the paper, the NTIR and DC image sets will correspond to domain A and domain B, respectively. Taking the translation from domain A to domain B as an example, the input image pair of domain A and B is denoted as  $\{x_a, x_b\}$ , the generator  $G_{AB}$  contains an encoder of domain A and a decoder of domain B, and the discriminator  $D_B$  aims to distinguish the real data  $x_b$  from the translated data  $G_{AB}(x_a)$ . Similarly, the inverse mapping includes the generator  $G_{BA}$  and the discriminator  $D_A$ .

### B. Revisiting and Improving ToDayGAN Model

1) *Revisiting ToDayGAN Model:* ToDayGAN modifies the ComboGAN model [37] to improve the performance of retrieval-based localization tasks. It has two paired generator-discriminator modules, which are capable of learning mappings between nighttime and daytime visible images. The generator networks are the same as the networks used in CycleGAN, whereas each discriminator network contains three copies to focus on different aspects (i.e., texture, color, and gradients) of the input. Similarly to CycleGAN, ToDayGAN model is supervised by two losses, i.e., adversarial loss and cycle-consistency loss  $\mathcal{L}_{cyc}^{ori}$ . Differently, it chooses Relativistic Loss [38] adapted to least-squares GAN loss as adversarial loss  $\mathcal{L}_{adv}$ .

2) *Improving ToDayGAN Model for TIR Image Colorization:* There are three distinct limitations when directly applying ToDayGAN to the NTIR2DC task. First of all, the translated results usually exhibit color dot artifacts, which degrade the naturalness of the image. Secondly, we observe

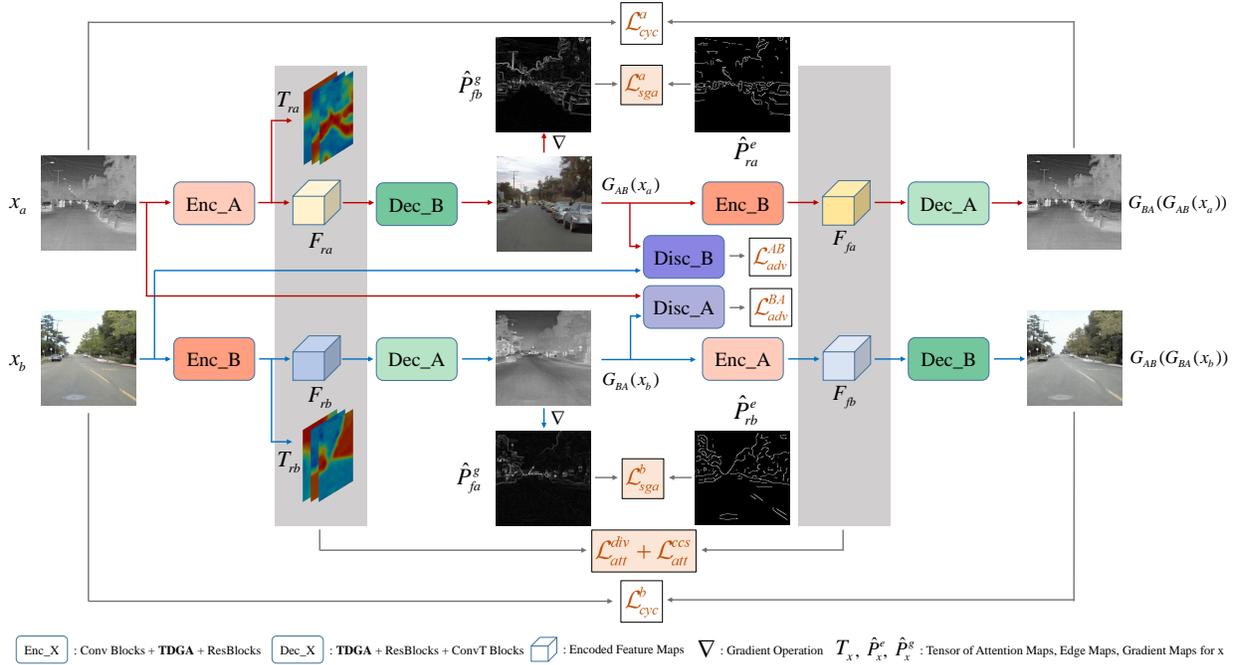


Fig. 2. The overall architecture of the proposed method.  $x_a$  and  $x_b$  respectively denote random images from nighttime TIR domain A and daytime visible domain B. The gray arrows indicate the composition of the loss function, and the red and blue arrows correspond to the forward computation of domain A and domain B, respectively. The abbreviations for Encoder, Decoder and Discriminator are Enc, Dec and Disc, respectively. We first propose a TDGA module and an elaborate attentional loss ( $\mathcal{L}_{att}^{div} + \mathcal{L}_{att}^{ccs}$ ) to reduce the semantic encoding ambiguity during translation. Then, a structured gradient alignment loss ( $\mathcal{L}_{sga}^a + \mathcal{L}_{sga}^b$ ) is introduced to encourage geometric consistency between the translated images and input images.

that sometimes the reconstructed TIR image is different from the input TIR image, although the cycle-consistency loss value is still small. Furthermore, the training process is unstable.

To avoid the aforementioned drawbacks, we adjust some designs of the ToDayGAN model. Referring to StyleGAN2 [39], droplet artifacts may be caused by the generator’s desire to use instance normalization to achieve scale-specific controls. But unlike StyleGAN2, which uses the weight demodulation module to replace instance normalization, we design a novel alternative that removes the artifacts while bringing fewer changes to the model structure. The main idea is replacing the last two instance normalization layers of the decoder with two group normalization [40] layers, and then introducing the total variation [41] loss  $\mathcal{L}_{tv}$  to punish the noise in the translated image. The problem of inaccurate representation of cycle-consistency loss motivates us to think about the expressive ability of loss function. We notice that the temperature of some background categories (e.g., tree, building and road) in the night environment is low and the difference among them is small, resulting in low contrast and blurred boundaries among local areas of the TIR image. We speculate that this may be the reason why the cycle-consistency loss based on the L1-norm is not sensitive to the structural differences between TIR images. Consequently, we additionally introduce SSIM [42] loss  $\mathcal{L}_{ssim}$  to supplement the evaluation of structural differences between images<sup>1</sup>. Then the improved cycle-consistency loss can be

expressed as:

$$\mathcal{L}_{cyc} = \lambda_{cyc} \mathcal{L}_{cyc}^{ori} + \lambda_{ssim} \mathcal{L}_{ssim}, \quad (1)$$

where  $\lambda_{cyc}$  and  $\lambda_{ssim}$  are loss weights. For the problem of unstable training, similar to [25], [26], we add spectral normalization [43] after each convolutional layer of the discriminator. Finally, combining the above three adjustments, we can obtain the variant model ToDayGAN-TIR, which serves as a starting point for the designed PearlGAN.

### C. Top-Down Guided Attention Module

The ToDayGAN-TIR model still does not solve the local encoding ambiguity problem of TIR images. In the absence of a wide range of contextual information, understanding local areas of TIR images is more challenging than visible images, such as recognizing the wheel of a car. We speculate that this may be the reason for the confusing encoding of local features of TIR images by existing image translation models. Inspired by the attentional feedback mechanism in biological vision [11], we design the TDGA module and the customized attentional loss, which not only gradually uses a wide range of contextual information to assist the semantic perception of complex regions, but also reduces the spatial entanglement of features.

An illustration of the TDGA module is provided in Fig. 3. For a given input feature map  $F \in \mathbb{R}^{c \times h \times w}$ , we first combine a convolutional layer to separate features of different scales to obtain four sets of features  $\{F_1, F_2, F_3, F_4\}$ , and the channel number of each set of features is  $\frac{c}{4}$ . Then, a statistical

<sup>1</sup>We use the implementation provided by <https://www.cnpython.com/pypi/pytorch-msssim>.

feature pyramid  $F_{d_s} \in \mathbb{R}^{\frac{h}{2^s} \times \frac{w}{2^s}}$  is obtained by using four scales of stacked  $2 \times$  average pooling operation to represent the statistical features of different receptive field ranges, where  $s \in \{1, 2, 3, 4\}$  indexes the pyramid scale. Next, the features of the large receptive field progressively predict the attention features of small scales, and realize the attention estimation and feature extraction from coarse to fine. Specifically, we apply a  $3 \times 3$  convolution and Sigmoid activation function to  $F$ , which produces a 2D spatial attention map  $A_F \in \mathbb{R}^{h \times w}$ :

$$A_F = \text{Att}(F) = \sigma(\text{conv}(F; \hat{\theta})), \quad (2)$$

where  $\sigma(\cdot)$  denotes the Sigmoid activation function, and  $\text{conv}(\cdot; \hat{\theta})$  denotes a convolutional layer with the parameter  $\hat{\theta}$ . In order to use a wide range of contextual information to gradually predict the attention map, we first utilize the smallest resolution feature map  $F_{d_4}$  to infer an attention map  $A_{d_4}$  by Eq. (2). The attention cue mined at this scale is applied over the feature of the next scale, which guides the features of the next scale to pay attention to contextual signal and predicts the attention map more completely. Then, the attention map at the second feature resolution (i.e.,  $8 \times$  down-sampling) is generated by:

$$A_{d_3} = \text{Att}(F_{d_3} + F_{d_3} \odot (A_{d_4} \uparrow_2)), \quad (3)$$

where  $\odot$  denotes element-wise multiplication with channel-wise broadcasting, and  $\uparrow_2$  means the  $2 \times$  spatial up-sampling operation. By analogy, we can infer the other attention maps  $A_{d_2}$  and  $A_{d_1}$  of different scales by Eq. (3). Through the cascaded architecture, higher-resolution features can focus on the areas with more integral semantics under the guidance of the lower-resolution attention cues, which provide a wide range of contextual information. To preserve the hierarchical attention features, the attention cue mined at each scale is applied over the corresponding feature of the original spatial resolution, which are further merged in a cascading manner. The output feature maps are generated by:

$$F' = \text{concat}((F_4 + F_4 \odot (A_{d_4} \uparrow_{16})), (F_3 + F_3 \odot (A_{d_3} \uparrow_8)), (F_2 + F_2 \odot (A_{d_2} \uparrow_4)), (F_1 + F_1 \odot (A_{d_1} \uparrow_2))), \quad (4)$$

where  $\text{concat}(\cdot)$  represents feature channel concatenation. Since the middle layer of the network may need enough contextual information to realize the transition from low-level features to high-level semantics, we simply put the TDGA module before the residual module in the encoder and decoder to assist with feature disentanglement.

#### D. Attentional Loss

To further ensure that the TDGA module can produce a hierarchical attention distribution, we introduce an elaborate attentional loss. Attentional loss is composed of AD loss and ACCS loss, which encourage hierarchical coding of features at the spatial distribution level and feature level, respectively. Since the global attention can further benefit robust perception of the whole scene with stronger generalization ability, the designed attentional loss is only applied to the three finer scales (i.e.,  $8 \times, 4 \times, 2 \times$  down-sampling). In other words, the number of attention scales applied to attentional loss, denoted

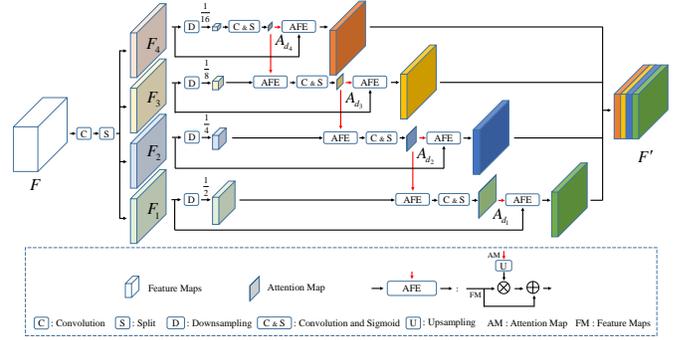


Fig. 3. Illustration of the proposed TDGA module. AFE indicates attention-directed feature enhancement.

as  $n_a$ , is set to three throughout the rest of the paper. We first restore attention maps of different scales to the original spatial resolution (i.e.,  $h \times w$ ) by an up-sampling operation, and then the attention maps are further concatenated to obtain the cascaded attention tensor  $T \in \mathbb{R}^{n_a \times h \times w}$ :

$$T = \text{concat}(A_{d_3} \uparrow_8, A_{d_2} \uparrow_4, A_{d_1} \uparrow_2). \quad (5)$$

The details of the introduced attentional loss are described below.

1) *Attentional Diversity Loss*: AD loss enforces the mutual exclusion and completeness of attention at multiple scales in space, and the concrete form is designed as follows:

$$\mathcal{L}_{att}^{div} = \alpha \times \frac{1}{h \times w} \times \sum_{i,j} [(1 - \max_{k \in S} T_{k,i,j}) + \beta \times (\sum_k T_{k,i,j} - 1)^2], \quad (6)$$

where  $\alpha$  and  $\beta$  are tuning coefficients to enforce the range of loss being 0 to 1,  $S = \{1, \dots, n_a\}$  represents the set for the scales of attention maps, and  $T_{k,i,j}$  denotes the attention weight located at  $(i, j)$  for the  $k_{th}$  channel of the attention tensor. In the square brackets, the two parts before and after the plus sign encourage the maximum value and sum value of the multi-scale attention maps at the same spatial position to be as close to one as possible. In this way, the attention of different scales will focus on different spatial regions, which benefits feature disentanglement. To normalize the output range to  $(0, 1)$ , the coefficients  $\alpha$  and  $\beta$  are set to  $\frac{1}{2}$  and  $\frac{1}{4}$ , respectively.

2) *Attentional Cross-Domain Conditional Similarity Loss*: Although AD loss encourages the diversified spatial distribution of attention, a lack of semantic level constraints might lead to irregular distribution of attention and hinder hierarchical coding of features. Hence, we introduce ACCS loss not only to encourage the compactness of the attentional features, but also to encourage the conditional similarity of the same-scale attentional features across the domains, which will be beneficial to generate a synthesized image with richer details and more natural texture. We first define the encoding features (i.e., the output of the encoder) of real NTIR, real DC, fake NTIR and fake DC images as  $F_{ra}$ ,  $F_{rb}$ ,  $F_{fa}$  and  $F_{fb}$ , respectively, and all of them have the same dimensions as  $F$ . Similar to the cascaded attention tensor  $T$ , the attention tensors of real NTIR and real DC images can be denoted as  $T_{ra}$  and  $T_{rb}$  through Eq. (5), and the dimensions of the two are the

same as T. We define the attention feature as the weighted sum of the attention map and the coding feature in the spatial dimension and normalize it with its L2 norm. For example, we suppose the attention feature of a fake DC image is denoted as  $V_{fbra} \in \mathbb{R}^{n_a \times c}$ , which is calculated by using  $F_{fb}$  and  $T_{ra}$ , and its  $k_{th}$  scale component  $V_{fbra}^k \in \mathbb{R}^{1 \times c}$  is defined as:

$$V_{fbra}^k = \frac{\tilde{V}_{fbra}^k}{\left\| \tilde{V}_{fbra}^k \right\|_2}, \quad (7)$$

where  $\|\cdot\|_2$  represents the L2 norm, and the unnormalized feature  $\tilde{V}_{fbra}^k$  can be formulated as:

$$\tilde{V}_{fbra}^k = \frac{GAP(F_{fb} \odot T_{ra}^k)}{GAP(T_{ra}^k)}, \quad (8)$$

where  $GAP(\cdot)$  denotes the globally averaged pooling operation, and  $T_{ra}^k$  represents the  $k_{th}$  channel of attention tensor  $T_{ra}$ . Then we can obtain different scale components of the attention feature (i.e.,  $V_{rara}$ ,  $V_{rbrb}$ ,  $V_{farb}$  and  $V_{fbra}$ ) for real-fake image pairs by utilizing Eq. (7). Next, we can calculate the similarity between attention features with cosine distance. For example, it is assumed that the similarity of attention features in the visible spectral domain can be defined as  $Q_b \in \mathbb{R}^{n_a \times n_a}$ , which is given by:

$$Q_b = Mm(V_{rbrb}, (V_{fbra})^T) = V_{rbrb} \otimes (V_{fbra})^T, \quad (9)$$

where  $\otimes$  is the matrix multiplication, and  $(\cdot)^T$  represents the transpose of the matrix. Based on the observation of the regularity of the spatial distribution of semantic categories in traffic scene images, for NTIR-DC image pairs with similar semantic distribution, we expect that the similarity of their attention features at the same scales should be greater than that of cross scales. Let  $diag(Q_b) \in \mathbb{R}^{n_a \times 1}$  be the diagonal elements of the similarity matrix  $Q_b$ , where the  $k_{th}$  element represents the same-scale similarity of the  $k_{th}$  scale, and let  $M(Q_b)$  denote the row-wise maximum of similarity matrix  $Q_b$ . Then we formulate ACCS loss function in the visible spectral domain as:

$$\mathcal{L}_{accs}^b = \frac{W_Q \otimes [M(Q_b) - diag(Q_b)]}{\sum_{k \in S} (W_Q)_k} + Dis(V_{rbrb}) + Dis(V_{fbra}), \quad (10)$$

where  $W_Q \in \mathbb{R}^{1 \times n_a}$  is an indicator vector representing the confidence of different scales of cross-domain attention, and  $Dis(\cdot)$  is a distance function to encourage a compact distribution of attention features. Concretely, the  $k_{th}$  element of  $W_Q$  is given by:

$$(W_Q)_k = \min(\max(T_{ra}^k), \max(T_{rb}^k)). \quad (11)$$

The distance function  $Dis(\cdot)$  expects the feature distance between scales to be as large as possible to ensure a compact distribution of features within the scale. For example, let  $Q_{rbrb} \in \mathbb{R}^{n_a \times n_a}$  denote the matrix product between  $V_{rbrb}$  and its transpose. Then the distance function for attention feature  $V_{rbrb}$  can be formulated as:

$$Dis(V_{rbrb}) = \max\left(\frac{\sum_{i \neq j, i, j \in S} (Q_{rbrb})_{ij}}{n_a(n_a - 1)}, 0\right), \quad (12)$$

where the numerator represents the sum of the elements on the non-main diagonal, the denominator represents the number of the non-main diagonal elements, and the the purpose of non-linear function  $\max(x, 0)$  is to avoid negative values. Similarly, the ACCS loss in the TIR spectral domain  $\mathcal{L}_{accs}^a$  can be obtained through Eq. (10). Finally, the complete ACCS loss function is defined as:

$$\mathcal{L}_{att}^{ccs} = \mathcal{L}_{accs}^a + \mathcal{L}_{accs}^b, \quad (13)$$

### E. Structured Gradient Alignment Loss

While the TDGA module and attentional loss can improve the stability of the model for TIR image encoding, the problem of edge distortion in the image translation process is still under-resolved. We observe that edge regions in TIR images usually appear at the junction of two regions with temperature difference, and the gradient at its corresponding position in the visible image is likely to be larger than the average of their neighboring gradients. In light of that, we propose a SGA loss to encourage the ratio of the gradient of the resulting image at the edge position to the maximum value of its neighborhood to be greater than a given threshold, while ignoring the constraints on non-edge regions.

Different from the previous methods [18], [23], the proposed SGA loss not only avoids the need for a pre-trained edge detection network, but also explicitly restricts the gradient structure of the translated image. Specifically, taking the TIR spectral domain as an example, we first use the offline edge detection method MCI [44] to predict the edge map of the TIR image in the training set. Then, we let the size of the image block be  $l_p \times l_p$  and use the method of adaptive average pooling to randomly select an image patch  $P_{ra}^e$  with edge pixels in the TIR edge map. The corresponding image patch in the gradient map of the fake DC image is named  $P_{fb}^g$ . Both  $P_{ra}^e$  and  $P_{fb}^g$  are normalized to the range  $[0, 1]$  by dividing the maximum of the corresponding image patch. Next, the SGA loss of the TIR domain can be defined as:

$$\mathcal{L}_{sga}^a = \frac{\sum_{i=1}^{l_p} \sum_{j=1}^{l_p} \max((\eta \times P_{ra}^e - P_{fb}^g)_{ij}, 0)}{\sum_{i=1}^{l_p} \sum_{j=1}^{l_p} P_{ra}^e}, \quad (14)$$

where the parameter  $\eta$  is a threshold that controls edge sharpness. Due to the variability of the brightness range of TIR images from different datasets, we set the value of parameter  $\eta$  to be related to the maximum intensity value of all infrared images in the dataset, which is denoted as  $I_{max}$  and given by:

$$\eta = 0.8 \times \frac{I_{max}}{255}, \quad (15)$$

where 0.8 is an empirical threshold representing the minimum gradient ratio to get a clear edge after normalization, and 255 represents the maximum value of an 8-bit image. By analogy, we can obtain the SGA loss of the visible image domain  $\mathcal{L}_{sga}^b$  by Eq. (14). Consequently, the SGA loss for two domains can be formulated as:

$$\mathcal{L}_{sga} = \mathcal{L}_{sga}^a + \mathcal{L}_{sga}^b, \quad (16)$$

Therefore, the proposed SGA loss can punish the result of inconspicuous or disappearing edges to maintain the consistency of edges between two domains.

#### F. Objective Function

In summary, the full objective function of PearlGAN can be written as:

$$\mathcal{L}_{all} = \mathcal{L}_{adv} + (\lambda_{cyc}\mathcal{L}_{cyc}^{ori} + \lambda_{ssim}\mathcal{L}_{ssim}) + \lambda_{tv}\mathcal{L}_{tv} + \lambda_{att}(\mathcal{L}_{att}^{div} + \mathcal{L}_{att}^{ccs}) + \lambda_{sga}\mathcal{L}_{sga}, \quad (17)$$

where  $\lambda_{cyc}$ ,  $\lambda_{ssim}$ ,  $\lambda_{tv}$ ,  $\lambda_{att}$  and  $\lambda_{sga}$  are loss weights. Referring to CycleGAN [7],  $\lambda_{cyc}$  is set to 10, and we empirically set the value of  $\lambda_{tv}$  to be half of  $\lambda_{cyc}$  to reduce color dot artifacts. Without loss of generality, the values of  $\lambda_{ssim}$  and  $\lambda_{att}$  are both set to one. However, the naturalness of the translated image is sensitive to the parameter  $\lambda_{sga}$ . If  $\lambda_{sga}$  were to be set to one, the generated image would have sharp edges but the naturalness would be heavily degraded. Hence, we simply set the value of  $\lambda_{sga}$  to 0.5 without more trials.

### IV. EXPERIMENTS

In this section, we first introduce the datasets and evaluation metrics for the NTIR2DC task. Subsequently, we explain the annotation of the test images and the implementation details of the model. Then, experimental results on FLIR and KAIST datasets are presented. Next, we perform an ablation analysis to verify the validity of the proposed module and loss functions. At last, some discussion of the experimental results is presented.

#### A. Datasets and Evaluation Metrics

1) *Datasets*: Experiments are conducted on the FLIR and KAIST datasets to demonstrate the effectiveness of our PearlGAN. The FLIR Thermal Starter Dataset [13] provides an annotated TIR image set and non-annotated RGB image set for training and validating object detection models. According to the lighting conditions of RGB images, we finally obtain 5447 DC images and 2899 NTIR images for training, and 490 NTIR images in the validation set for model evaluation.

The KAIST Multispectral Pedestrian Detection Benchmark [14] provides somewhat aligned color and thermal image pairs captured in day and night. Due to the low brightness of some DC images, as shown in Fig. 4, we enhance all the training DC images in KAIST using the SRLIE [45] method to improve the image quality. We finally obtain 1674 DC images and 1359 NTIR images as training set samples, and 500 NTIR images in the test set for evaluating semantic and edge consistency<sup>2</sup>. The sample size of the test set used for pedestrian detection experiments is 611.

Due to differences in sensor imaging resolution and processing of image registration, black filled areas appear on both sides of some images in the FLIR and KAIST datasets. In order to remove irrelevant regions, we first resize all training set images to  $500 \times 400$  resolution, and then perform center crop to get  $360 \times 288$  resolution input images.

<sup>2</sup>See <https://github.com/FuyaLuo/PearlGAN/> for specific sample selection.

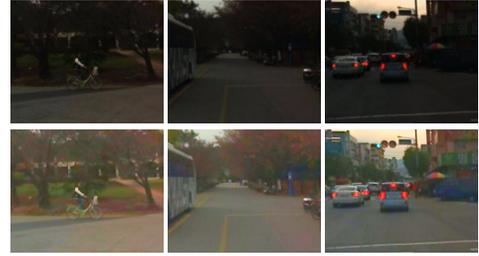


Fig. 4. KAIST image enhancement results achieved by SRLIE [45]. The first and second rows are the original image and the enhanced image, respectively.

2) *Evaluation Metrics*: An ideal NTIR image colorization method should preserve the image content at all levels, i.e., from scene level layout to fine-grained edges. Therefore, inspired by [28], we perform three vision tasks to evaluate the coloring performance of NTIR images, including semantic segmentation, object detection, and edge detection.

Intersection-over-Union (IoU) [46], which is the ratio of the intersection of the predicted segmentation map with the ground truth to their union, is a widely used metric for semantic segmentation. We take the mean IoU (mIoU) over all categories to indicate the overall performance of the model.

For object detection, average precision (AP) [47] is defined as the average detection precision under different recalls. In the case of multiple categories, the mean AP (mAP) averaged over all object categories is typically used as the final metric of performance.

Since the model may generate more details to obtain plausible DC images, it is sub-optimal to use traditional metrics (e.g., F-score, which considers both precision and recall) to evaluate the edge consistency of the NTIR2DC task. Therefore, we propose a new metric called APCE to evaluate the preservation of edges in the source NTIR image. APCE is the average precision of Canny edges under different threshold conditions. Given the  $j$ th high threshold  $\mu_j$ , the Canny edge maps of the  $i$ th test image and its corresponding output image are denoted as  $X_i(\mu_j)$  and  $Y_i(\mu_j)$ , respectively. Then, we define the APCE of the entire test set as:

$$APCE = \frac{1}{n_i} \frac{1}{n_j} \sum_{i=1}^{n_i} \sum_{j=1}^{n_j} \frac{X_i(\mu_j) \cap Y_i(\mu_j)}{X_i(\mu_j)}, \quad (18)$$

where  $n_i$  and  $n_j$  represent the total number of images and thresholds, respectively.

#### B. Image Annotation and Implementation Details

1) *Image Annotation*: To evaluate the semantic preservation performance of the NTIR2DC task, a subset of FLIR and KAIST datasets are selected and annotated with their pixel-level category labels. Due to the low contrast and ambiguous boundaries of NTIR images, as shown in Fig. 5, we define nine categories and use the *LabelMe* toolbox<sup>3</sup> to annotate only their identified corresponding regions. The labeled categories are road, building, traffic sign, sky, people, car, truck, bus and motorcycle.

<sup>3</sup><https://github.com/CSAILVision/LabelMeAnnotationTool>.

TABLE I  
SEMANTIC SEGMENTATION RESULTS OF THE SYNTHESIZED IMAGES OBTAINED BY DIFFERENT TRANSLATION METHODS ON FLIR DATASET. ALL NUMBERS ARE IN %

	Road	Building	Traffic sign	Sky	Person	Car	Truck	Bus	Motorcycle	mIoU
Reference NVC images	95.5	62.8	<b>6.4</b>	63.4	40.7	51.4	0.0	<b>0.4</b>	0.0	35.6
CycleGAN [7]	95.6	39.1	0.0	90.8	60.8	78.0	0.0	0.0	0.0	40.5
UNIT [24]	96.2	60.3	0.2	92.1	64.5	71.5	0.0	0.0	<b>14.6</b>	44.4
MUNIT [8]	96.0	27.5	0.0	92.8	49.6	64.3	0.0	0.0	0.0	36.7
ToDayGAN [9]	95.7	47.2	0.0	85.3	56.8	75.4	0.0	0.0	0.0	40.0
UGATIT [25]	94.3	18.4	0.0	89.0	23.7	59.0	0.0	0.0	0.0	31.6
DRIT++ [27]	97.3	29.4	0.0	78.4	28.0	78.9	0.0	0.0	0.0	34.7
ForkGAN [28]	94.4	55.1	0.0	90.7	60.5	76.1	0.0	0.0	0.0	41.9
Proposed	<b>97.7</b>	<b>73.1</b>	0.0	<b>93.4</b>	<b>73.2</b>	<b>82.7</b>	<b>0.1</b>	0.0	0.0	<b>46.7</b>

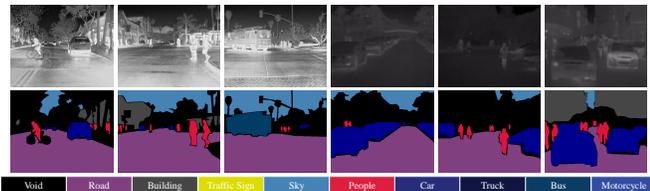


Fig. 5. Examples of NTIR images from the FLIR (the first three columns) and KAIST (the last three columns) datasets and our annotated segmentation masks.

2) *Experiment Settings and Implementation Details:* Since there is little available source code for the NTIR2DC task, we compare the proposed PearlGAN with other state-of-the-art I2I translation methods, including CycleGAN [7], UNIT [24], MUNIT [8], DRIT++ [27] and UGATIT [25], and low-light enhancement methods such as ToDayGAN [9] and ForkGAN [28]. We follow the instructions of these methods in order to make a fair setting for comparison.

The proposed PearlGAN is implemented using PyTorch. We train the models using the Adam optimizer [48] with  $(\beta_1, \beta_2) = (0.5, 0.999)$  on NVIDIA RTX 3090 GPUs. The batch size is set to one for all experiments. The learning rate begins at 0.0002, is constant for the first half of training and decreases linearly to zero during the second half of training. The total number of training epochs for the FLIR and KAIST datasets are 80 and 120, respectively. Due to the need to focus on optimizing the cycle-consistency loss early in the training, the SSIM loss and ACCS loss are set to start working after about 50K iterations (i.e., around the 10th and 30th epochs of the FLIR and KAIST dataset training, respectively) in order to stabilize the training process. In Eq. (14), the side length  $l_p$  of the image patch is set to 32, and the gradient ratio thresholds  $\eta$  for the FLIR and KAIST datasets are obtained as 0.8 and 0.44, respectively, from Eq. (15). For data augmentation, we flip the images horizontally with a probability of 0.5, and randomly crop them to  $256 \times 256$ .

Due to the lack of pixel-level semantic labels in the FLIR and KAIST datasets, we perform domain adaptive semantic segmentation using the Cityscape dataset [49] and the advanced scene adaptation method [50] on real DC images from both datasets, and the final trained models are used for evaluation of the semantic segmentation performance of the translated images. To measure the generalizability of the features of the objects in the synthesized images, we utilize

the YOLOv4 [51] model pre-trained on the MS COCO dataset [52] as an object detection evaluation method. For the APCE evaluation metric, the high threshold of Canny edge detection<sup>4</sup> ranges from 0.01 to 0.99 with an interval of 0.01, and the low threshold is taken as half of the high threshold.

### C. Experiments on FLIR Dataset

For a comprehensive evaluation of our method, we conduct three experiments on FLIR datasets, including semantic segmentation, object detection and edge preservation.

1) *Semantic Segmentation:* The translated results and corresponding segmentation outputs from using various methods on the FLIR dataset are presented in Fig. 6. Column (a) lists the reference nighttime visible color (NVC) image and its semantic segmentation output. Although NVC images can provide somewhat reasonable layout cues, small cars in the distance with strong lighting are extremely challenging for the segmentation model to identify. In contrast, the translated images obtained by most translation methods can facilitate the localization of distant small cars. However, as shown by the white dashed box in the segmentation masks, CycleGAN, MUNIT, UGATIT, DRIT++ and ForkGAN fail to generate plausible pedestrians, and the translated images of UNIT and TDG do not guarantee complete pedestrian mask prediction. In contrast, our method has a stronger ability to preserve scene layout and generate plausible textures (e.g., buildings and cars). Table I reports a quantitative comparison of the semantic consistency performance on the FLIR dataset. The proposed method achieves the highest mIoU (46.7%) among all the methods. In addition, our approach achieves the best results in the translation of semantic regions in five categories: road, building, sky, person, and car. Note that due to the small number of samples containing trucks, buses and motorcycles in the training set, as well as the small area occupied by traffic signs, all methods have poor translation results for these four categories.

2) *Object Detection:* Fig. 7 shows the qualitative translation and detection result comparisons, wherein the second row is the zoomed in image of the corresponding area. As shown in the red dashed box, UNIT, MUNIT, DRIT and UGATIT cannot generate plausible cars, while CycleGAN, ToDayGAN and ForkGAN fail to maintain distant cars. In contrast, our

<sup>4</sup>We use the implementation provided by <https://www.mathworks.com/help/images/ref/edge.html?lang=en>.

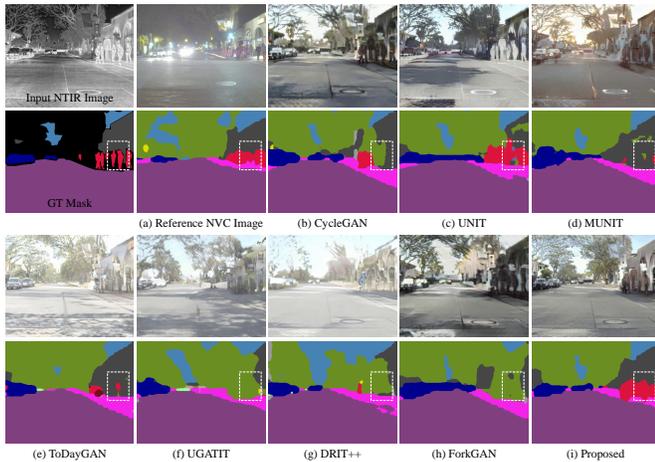


Fig. 6. The visual translation (the first row) and segmentation performance (the second row) comparison of different methods on the FLIR dataset. Please zoom in to check more details on the content and quality. The area covered by the white dotted box is worth attention.



Fig. 7. Visual comparison of detection results on the FLIR dataset by YOLOv4 model [51]. The parts covered by red and green dashed boxes show the enlarged cropped region in the corresponding image. Colors in the detection results that do not intersect with GT represent undefined categories of FLIR dataset identified by the detector.

method not only generates quite realistic cars, but also significantly outperforms other methods in terms of small object preservation. Furthermore, as shown in the green dashed box, all the I2I image translation methods are unable to generate reasonable pedestrian features to convince the general object detector except ours. In addition, our translation results can help more objects to be detected compared with the original NVC images. Since only three of the annotation categories (i.e., person, bicycle and car) provided by the FLIR dataset exist in the validation set, a quantitative comparison of the detection performance on the FLIR dataset is listed in Table II. As shown, the detection performance on the translated images obtained by the proposed method substantially outperforms other methods, and our mAP result is almost twice as good as the second ranked method (i.e., 50.8 vs 24.8), which illustrates the superiority of our method in terms of object preservation.

3) *Edge Preservation*: In Fig. 8, we qualitatively compare the edge preservation of different image translation methods, and the second row shows the zoomed-in patches of the corresponding area fused with the edges of the original NTIR

TABLE II  
OBJECT DETECTION RESULTS OF THE SYNTHESIZED IMAGES OBTAINED BY DIFFERENT TRANSLATION METHODS ON FLIR DATASET, COMPUTED AT A SINGLE IOU OF 0.50. ALL NUMBERS ARE IN %

	Person	Bicycle	Car	mAP
Reference NVC images	9.8	2.6	11.5	8.0
CycleGAN [7]	17.8	1.9	37.2	19.0
UNIT [24]	16.3	9.5	18.3	14.7
MUNIT [8]	29.5	2.3	42.5	24.8
ToDayGAN [9]	19.0	1.5	53.3	24.6
UGATIT [25]	1.9	0.0	14.6	5.5
DRIT++ [27]	16.5	2.2	46.0	21.6
ForkGAN [28]	25.9	2.3	32.5	20.2
Proposed	<b>54.0</b>	<b>23.0</b>	<b>75.5</b>	<b>50.8</b>

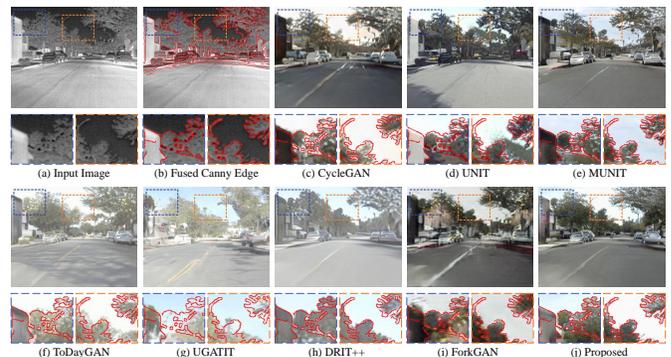


Fig. 8. Visual comparison of geometric consistency results on FLIR dataset. The second row shows the enlarged results of the corresponding regions after fusion with the edges of the input image. Column (b) is the result of fusing the input image with its Canny edges.

image. All compared I2I translation methods commonly have edges that expand outward (e.g., the yellow dotted box area in columns d and f) or edges that contract inward (e.g., the blue dotted box area in columns d, e and g). In contrast, the translated images obtained by the proposed method better preserve the edge structure of the original images, as shown in the blue dashed box in column j. For the yellow dashed box region, although a portion of the branches generated by our method is beyond the edges of the original image (i.e., the red line), these branches are present in the original NTIR image but are not captured by the Canny edge detector with the current parameters. Therefore, we investigate the edge preservation performance under different Canny thresholds, and the results of each method are shown in Fig. 9(a). Considering all thresholds, our method consistently keeps higher performance in terms of edge consistency during translation.

#### D. Experiments on KAIST Dataset

To further verify the effectiveness and robustness of our proposed method, we conduct experiments on the challenging KAIST dataset. Due to the low resolution of early imaging equipment, both NTIR and DC images in the KAIST dataset are blurred and low contrast, which makes the NTIR2DC task even more difficult on this dataset.

1) *Semantic Segmentation*: Fig. 10 presents the translated results and the segmentation outputs of various methods on the KAIST dataset. As shown in the white dashed box in the

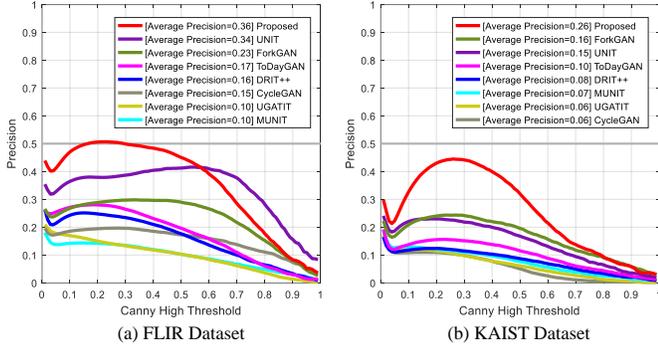


Fig. 9. APCE results of different translation methods on FLIR and KAIST datasets.

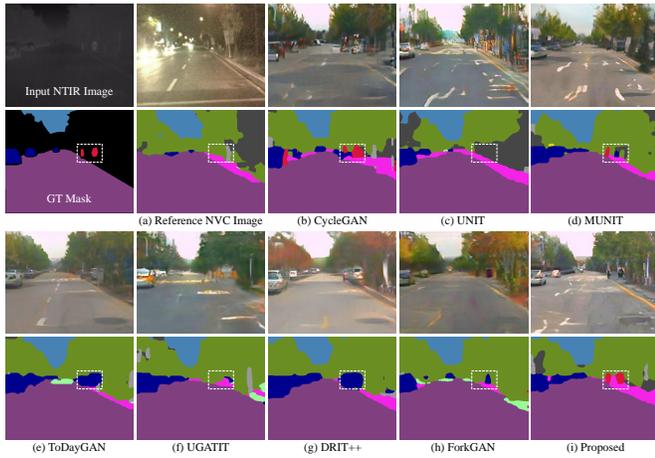


Fig. 10. The visual translation (the first row) and segmentation performance (the second row) comparison of different methods on KAIST dataset. The area covered by the white dotted box is worth attention.

second row, the segmentation model is unable to recognize the pedestrians on the roadside due to the low illumination in the NVC image. UNIT, ToDayGAN, UGATIT, DRIT++ and ForkGAN fail to capture the characteristics of pedestrians. Although CycleGAN and MUNIT can make the segmentation model perceive pedestrians, the generation of pedestrian regions is incomplete (e.g., MUNIT) or overfilled (e.g., CycleGAN). However, benefitting from the proposed TDGA module and SGA loss, our PearlGAN can better avoid the corruption of objects. The quantitative comparison is reported in Table III. After the alignment process with the NTIR images, the mIoU is only 33.8% when we directly perform the semantic segmentation on the real NVC images. Note that the proposed PearlGAN achieves the highest mIoU among all the methods, and a significant improvement (i.e., + 9.3%) compared with that on the real NVC images, which indicates that PearlGAN can facilitate nighttime scene perception while better preserving the scene layout. Similar to the experiments on the FLIR dataset, the poor performance in the translation of traffic signs and vehicles other than cars is due to the small number of samples available for learning.

2) *Pedestrian Detection*: Since the KAIST dataset only provides bounding box annotation for pedestrians, we investigate the pedestrian generation performance of different I2I

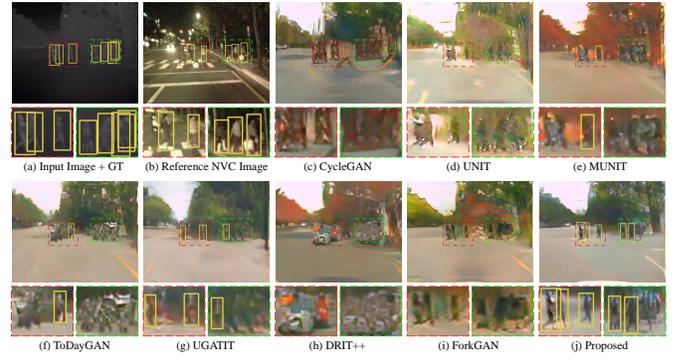


Fig. 11. Visual comparison of pedestrian detection results on KAIST dataset by YOLOv4 model [51]. The parts covered by red and green dashed boxes show the enlarged cropped regions in the corresponding image.

translation methods using YOLOv4. In Fig. 11, we qualitatively compare the different methods. Among all the compared methods, only MUNIT, ToDayGAN and UGATIT generate pedestrians accepted by the detector despite the lack of realistic structure, whereas the other methods fail completely. In contrast, as shown in the red dashed box in the figure, the proposed PearlGAN can better maintain the pedestrian structure and perform a more reasonable translation. The quantitative comparison is listed in Table IV. Due to the low quality of NTIR images and the small number of pedestrian samples in real DC images, the pedestrian detection performance on DC images synthesized by all translation methods is far inferior to that on the original NVC images. Nevertheless, the proposed method significantly outperforms the other compared translation methods in pedestrian transformation, and the mAP for pedestrian detection is twice as high as that of the second ranked method (i.e., 25.8 vs. 11.0).

3) *Edge Preservation*: The qualitative comparison of edge preservation on the KAIST dataset is shown in Fig. 12. For better visualization, column (b) shows the fusion results of the Canny edges of the original image and its enhanced image. We can observe from the figure that the compared translation methods all have edge shift problems, especially UNIT, MUNIT and DRIT++ (e.g., blue dashed box area). While the synthesized image obtained by the proposed method can better preserve the edge structure of the original image, as shown in column (j), the edges of the original NTIR image fit perfectly with the edges of the translated DC image. Considering different thresholds of Canny edges, the quantitative comparison of the edge consistency of various translation methods on the KAIST dataset is shown in Fig. 9(b). As shown, the edge consistency performance of the proposed method significantly surpasses other methods at almost all thresholds, which further illustrates the superiority of PearlGAN in edge structure preservation.

### E. Ablation Study

An ablation study is conducted on the FLIR dataset to evaluate the contribution of each component in PearlGAN. The results of the ablation study are presented in Table V, and an example of qualitative comparison of each proposed

TABLE III  
SEMANTIC SEGMENTATION RESULTS OF THE SYNTHESIZED IMAGES OBTAINED BY DIFFERENT TRANSLATION METHODS ON KAIST DATASET. ALL NUMBERS ARE IN %

	Road	Building	Traffic sign	Sky	Person	Car	Truck	Bus	Motorcycle	mIoU
Reference NVC images	92.2	71.8	0.0	66.3	15.7	57.7	0.0	0.2	0.0	33.8
CycleGAN [7]	88.0	48.7	0.0	78.6	15.0	49.2	0.0	0.0	0.0	31.1
UNIT [24]	94.1	<b>73.9</b>	3.1	86.6	36.0	67.7	0.0	0.0	1.6	40.3
MUNIT [8]	88.7	34.5	0.2	81.0	7.8	46.2	0.0	0.0	0.6	28.8
ToDayGAN [9]	93.3	63.2	2.3	87.7	20.4	58.3	0.0	0.0	0.0	36.1
UGATIT [25]	90.0	52.2	1.3	73.3	16.7	53.0	0.0	0.0	0.0	31.8
DRIT++ [27]	91.2	71.5	0.0	73.8	5.1	56.2	0.0	0.0	0.0	33.1
ForkGAN [28]	93.9	54.3	0.9	87.0	22.7	66.2	0.0	0.0	<b>2.9</b>	36.4
Proposed	<b>94.7</b>	72.2	<b>5.8</b>	<b>91.2</b>	<b>43.0</b>	<b>67.7</b>	0.0	<b>13.0</b>	0.0	<b>43.1</b>

TABLE IV  
PEDESTRIAN DETECTION RESULTS OF THE SYNTHESIZED IMAGES OBTAINED BY DIFFERENT TRANSLATION METHODS ON KAIST DATASET, COMPUTED AT A SINGLE IOU OF 0.50. ALL NUMBERS ARE IN %. TOP THREE RESULTS ARE MARKED IN RED, BLUE, AND GREEN.

	Precision	Recall	mAP
Reference NVC images	36.8	50.1	<b>44.2</b>
CycleGAN [7]	4.7	2.8	1.1
UNIT [24]	26.7	14.5	<b>11.0</b>
MUNIT [8]	2.1	1.6	0.3
ToDayGAN [9]	11.4	14.9	5.0
UGATIT [25]	13.3	7.6	3.2
DRIT++ [27]	7.9	4.1	1.2
ForkGAN [28]	33.9	4.6	4.9
Proposed	21.0	39.8	<b>25.8</b>



Fig. 12. Visual comparison of geometric consistency results on KAIST dataset. The second row shows the enlarged results of the corresponding regions after fusion with the edges of the input image. Column (b) is the result of blending the enhanced image of the input image with its Canny edges for better viewing.

component is shown in Fig. 13. We can find that the brightness of the synthesized image obtained from the baseline model is so high that it produces a hazy visual effect, and there is a degree of content distortion in the results. In addition, the output shows artifacts of multiple colored dots aggregated together (i.e., the purple box in the first row). Then, as shown in the third column of the figure, with the help of the existing modules, our redesigned ToDayGAN-TIR model not only eliminates artifacts but also improves image quality as well as reduces content distortion. The results in the table further illustrate the validity of the model adaptation. In the case of simply introducing the TDGA module without other

losses, although the model does not gain much in terms of content preservation, there is a significant improvement in edge consistency due to the coarse-to-fine feature encoding, which may help the model to better capture the edge information in NTIR images. In the next experiment, we study the influence of the attentional diversity loss. We find that the model gains a slight improvement in semantic preservation, when compared with the model using only the TDGA module. Then we investigate the effectiveness of the attentional cross-domain condition similarity loss. As shown in the red box in the sixth column of the figure, the translated image obtains a more realistic building area compared with the previous model. This illustrates that the introduced ACCS loss can reduce the feature entanglement in the model encoding process and achieve more reasonable translation. Finally, we explore the improvement from the proposed SGA loss. As shown in the figure, the edge structure of the translated image matches well with the original image after the SGA loss is introduced, which indicates that the proposed loss can effectively reduce geometric distortion. In addition, we find from the table that the reliable edge structure is beneficial for semantic preservation, which further demonstrates the superiority of the proposed architecture.

TABLE V  
QUANTITATIVE COMPARISONS FOR ABLATION STUDIES ON FLIR DATASET. "MA" MEANS THE MODEL ADAPTATION TO OBTAIN TODAYGAN-TIR.

Baseline	MA	TDGA	AD	ACCS	SGA	mIoU (%)	mAP (%)	APCE
✓						40.0	24.6	0.17
✓	✓					43.9	45.0	0.27
✓	✓	✓				44.3	46.2	0.34
✓	✓	✓	✓			44.8	48.1	0.34
✓	✓	✓	✓	✓		45.2	48.7	0.34
✓	✓	✓	✓	✓	✓	<b>46.7</b>	<b>50.8</b>	<b>0.36</b>

## F. Discussion

In this section, we first visualize the attention maps learned by the model, then discuss the FID results, and finally analyze the failure cases.

1) *Attention Map Visualization*: We visualize the attention maps of a pair of DC and NTIR images from the test set of the FLIR dataset to see whether they can learn spatially separated attention. The attention maps are presented in Fig. 14. As shown in the second column of the figure, the models all tend to pay attention to the top region of the image to

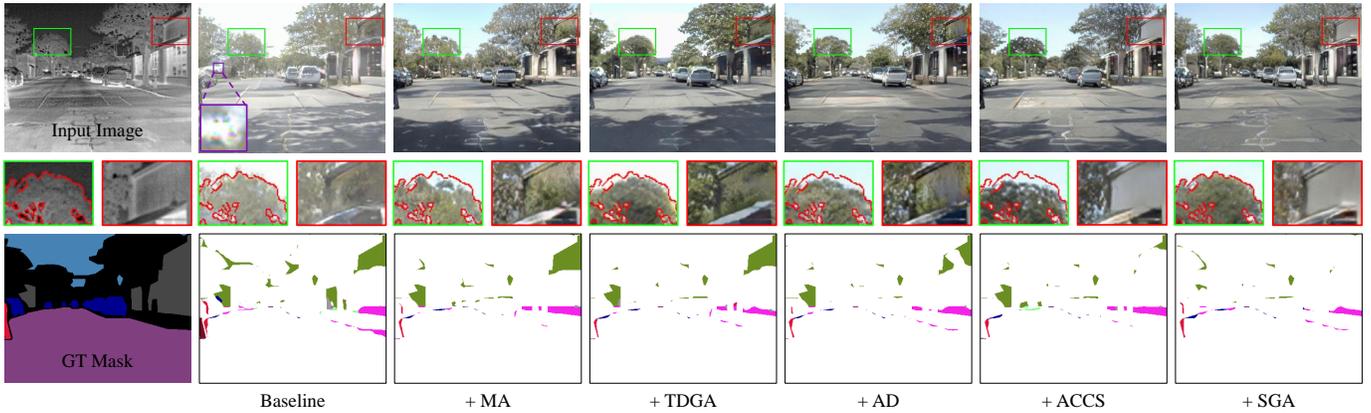


Fig. 13. Visual results of ablation study on FLIR dataset. The first row shows the translation results for different models. In the second row, the parts covered by green dashed boxes show the enlarged results of the corresponding regions after fusion with the edges of the input image, and the parts covered by red dashed boxes show the enlarged cropped regions in the corresponding image. The third row shows the error maps of the semantic segmentation results, where the white areas indicate the correct regions or unlabeled regions.

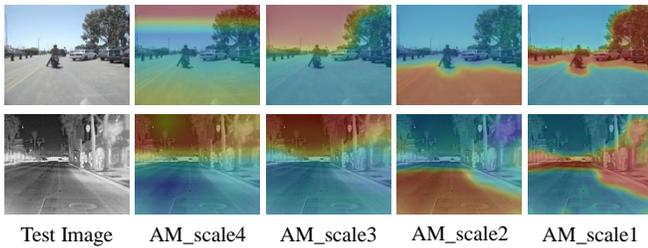


Fig. 14. Attention map visualization learned through TDGA module. The second to the fifth column indicate the attentional maps from the fourth to the first scale, respectively.

capture contextual information regardless of whether the image is a visible or infrared image. From the third to the first scale of the attentional map, the model tends to pay attention to the top, bottom and middle regions of the image, which correspond to the sky, road and object-related regions, respectively. This hierarchical feature encoding manner is beneficial for reducing the semantic entanglement in the neighboring space by using contextual information. The aforementioned experiments on FLIR and KAIST datasets demonstrate the effectiveness of this attention pattern.

2) *FID Results*: Due to its widespread adoption in I2I image translation performance evaluation, we report the FID scores of different methods for the NTIR2DC task in Table VI. As shown in the table, the FID score of the proposed method ranks in the middle of all compared methods, while MUNIT and ToDayGAN outperform PearlGAN on both datasets. However, the FID score is biased, as it is related to the number of samples [53]. Moreover, the FID score only considers the similarity of the feature distribution of all samples and fails to measure differences in content and geometry before and after translation [12]. As shown in Fig. 15, MUNIT and TDG translate people into vegetation (i.e., the first row) or road (i.e., the second row) in order to obtain more realistic texture information, which is unacceptable for the domain adaptation scenarios, requiring rigorous preservation of image content. In contrast, our approach encourages the model to maintain

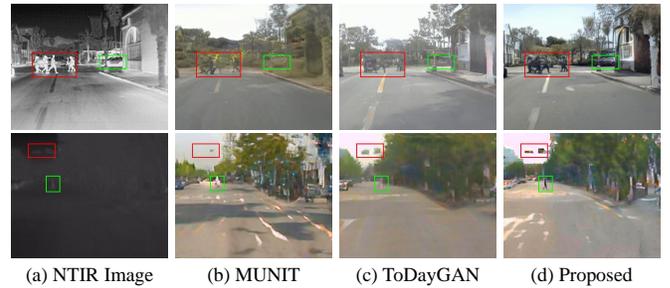


Fig. 15. Qualitative comparison of translation results. Areas covered by the red and green boxes are worth attention.

the geometry and reduce the semantic encoding ambiguity of the original image, and thus may generate local unnatural textures, leading to degradation of the FID score. Therefore, establishing a metric that simultaneously balances texture naturalness, content preservation and geometric consistency is an imperative research direction for future unpaired I2I translation work.

TABLE VI  
FID RESULTS ON FLIR AND KAIST DATASETS

	FLIR Dataset	KAIST Dataset
CycleGAN [7]	76.0	132.7
UNIT [24]	78.1	91.9
MUNIT [8]	<b>39.0</b>	98.3
ToDayGAN [9]	56.9	<b>90.7</b>
UGATIT [25]	69.8	99.8
DRIT++ [27]	58.9	105.9
ForkGAN [28]	99.8	175.3
Proposed	62.7	102.2

3) *Failure Cases*: Fig. 16 shows some failure cases of PearlGAN, the first two and last two example images are from the FLIR and KAIST datasets, respectively. As shown in the first and third columns of the figure, the proposed method fails to generate plausible buses due to few training samples containing this category. However, poor translation on small sample categories is a common defect among all comparison



Fig. 16. Visualization of failure cases. The first and second rows show the NTIR images and their translation results, respectively.

methods, as shown in Table I and Table III. A direct solution to this problem could be increasing the number of images in small sample categories by data augmentation. Although the proposed method can obtain natural background regions, reasonably translating objects in some complex and crowded scenes is still challenging, as shown in the second and fourth columns. To address this issue, more attempts should be made to develop better learning strategies for the understanding of scene layout in the NTIR images.

## V. CONCLUSION

In this paper, we propose a novel framework called PearlGAN to achieve TIR image colorization, which is beneficial to multiple vision tasks in nighttime driving scenes. Benefiting from the top-down guided attention structure and elaborated attentional loss, PearlGAN can learn hierarchical attention to reduce the spatial entanglement of features and better preserve semantic information. Moreover, we propose a structured gradient alignment loss to encourage geometric consistency between the translation result and the original image. In addition, we annotate a subset of FLIR and KAIST datasets with pixel-wise category labels to further catalyze research on colorization and semantic segmentation of NTIR images. Furthermore, we introduce a new evaluation metric to assess the edge consistency of the translation method. Comprehensive experiments demonstrate the superiority of PearlGAN for semantic preservation and edge consistency in the NTIR2DC task. In the future, designing a more reliable image translation model to maintain semantic consistency is a promising direction for our further research.

## REFERENCES

- [1] J. A. A. Cavanillas, "The role of color and false color in object recognition with degraded and non-degraded images," NAVAL POSTGRADUATE SCHOOL MONTEREY CA, Tech. Rep., 1999.
- [2] M. T. Sampson, "An assessment of the impact of fused monochrome and fused color night vision displays on reaction time and accuracy in target detection." NAVAL POSTGRADUATE SCHOOL MONTEREY CA, Tech. Rep., 1996.
- [3] U. Qayyum, Q. Ahsan, Z. Mahmood, and M. A. Chcmdary, "Thermal colorization using deep neural network," in *Proc. IBCAST*, 2018, pp. 325–329.
- [4] A. Berg, J. Ahlberg, and M. Felsberg, "Generating visible spectrum images from thermal infrared," in *Proc. CVPR Workshops*, 2018, pp. 1143–1152.
- [5] A. Nyberg, A. Eldesokey, D. Bergstrom, and D. Gustafsson, "Unpaired thermal to visible spectrum transfer using adversarial training," in *Proc. ECCV*, 2018, pp. 0–0.
- [6] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio, "Generative adversarial nets," in *Proc. NeurIPS*, 2014.
- [7] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proc. ICCV*, 2017, pp. 2223–2232.
- [8] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz, "Multimodal unsupervised image-to-image translation," in *Proc. ECCV*, 2018, pp. 172–189.
- [9] A. Anosheh, T. Sattler, R. Timofte, M. Pollefeys, and L. Van Gool, "Night-to-day image translation for retrieval-based localization," in *ICRA*, 2019, pp. 5958–5964.
- [10] Y. B. Saalmann, I. N. Pigarev, and T. R. Vidyasagar, "Neural mechanisms of visual attention: how top-down feedback highlights relevant locations," *Science*, vol. 316, no. 5831, pp. 1612–1615, 2007.
- [11] T. R. Vidyasagar, "A neuronal model of attentional spotlight: parietal guiding the temporal," *Brain Research Reviews*, vol. 30, no. 1, pp. 66–76, 1999.
- [12] 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," in *Proc. NeurIPS*, 2017, pp. 6629–6640.
- [13] F.A.Group, "Flir thermal dataset for algorithm training," <https://www.flir.co.uk/oem/adas/adas-dataset-form/>, May 2019.
- [14] S. Hwang, J. Park, N. Kim, Y. Choi, and I. So Kweon, "Multispectral pedestrian detection: Benchmark dataset and baseline," in *Proc. CVPR*, 2015, pp. 1037–1045.
- [15] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, 2016, pp. 770–778.
- [16] Y.-J. Cao, C. Lin, and Y.-J. Li, "Learning crisp boundaries using deep refinement network and adaptive weighting loss," *IEEE Trans. Multimedia*, vol. 23, pp. 761–771, 2021.
- [17] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proc. ICCV*, 2017, pp. 2961–2969.
- [18] T. Wang, T. Zhang, L. Liu, A. Wiliem, and B. Lovell, "Cannygan: Edge-preserving image translation with disentangled features," in *ICIP*, 2019, pp. 514–518.
- [19] R. Abbott, N. M. Robertson, J. Martinez del Rincon, and B. Connor, "Unsupervised object detection via lwir/rgb translation," in *Proc. CVPR Workshops*, 2020, pp. 90–91.
- [20] F. Almasri and O. Debeir, "Robust perceptual night vision in thermal colorization," *arXiv preprint arXiv:2003.02204*, 2020.
- [21] N. Bhat, N. Saggi, S. Kumar *et al.*, "Generating visible spectrum images from thermal infrared using conditional generative adversarial networks," in *ICCES*, 2020, pp. 1390–1394.
- [22] X. Kuang, J. Zhu, X. Sui, Y. Liu, C. Liu, Q. Chen, and G. Gu, "Thermal infrared colorization via conditional generative adversarial network," *Infrared Physics & Technology*, p. 103338, 2020.
- [23] T. Wang, T. Zhang, and B. C. Lovell, "Ebit: Weakly-supervised image translation with edge and boundary enhancement," *Pattern Recognition Letters*, vol. 138, pp. 534–539, 2020.
- [24] M.-Y. Liu, T. Breuel, and J. Kautz, "Unsupervised image-to-image translation networks," in *Proc. NeurIPS*, 2017.
- [25] J. Kim, M. Kim, H. Kang, and K. H. Lee, "U-gat-it: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation," in *Proc. ICLR*, 2019.
- [26] R. Chen, W. Huang, B. Huang, F. Sun, and B. Fang, "Reusing discriminators for encoding: Towards unsupervised image-to-image translation," in *Proc. CVPR*, 2020, pp. 8168–8177.
- [27] H.-Y. Lee, H.-Y. Tseng, Q. Mao, J.-B. Huang, Y.-D. Lu, M. Singh, and M.-H. Yang, "Drit++: Diverse image-to-image translation via disentangled representations," *International Journal of Computer Vision*, vol. 128, no. 10, pp. 2402–2417, 2020.
- [28] Z. Zheng, Y. Wu, X. Han, and J. Shi, "Forkgan: Seeing into the rainy night," in *Proc. ECCV*, 2020.
- [29] W. Wang, S. Zhao, J. Shen, S. C. Hoi, and A. Borji, "Salient object detection with pyramid attention and salient edges," in *Proc. CVPR*, 2019, pp. 1448–1457.
- [30] T. Zhao and X. Wu, "Pyramid feature attention network for saliency detection," in *Proc. CVPR*, 2019, pp. 3085–3094.
- [31] Q. Zhang, R. Cong, C. Li, M.-M. Cheng, Y. Fang, X. Cao, Y. Zhao, and S. Kwong, "Dense attention fluid network for salient object detection in optical remote sensing images," *IEEE Trans. Image Process.*, 2020.
- [32] Q. Wang, T. Wu, H. Zheng, and G. Guo, "Hierarchical pyramid diverse attention networks for face recognition," in *Proc. CVPR*, 2020, pp. 8326–8335.
- [33] D. Xie, C. Deng, H. Wang, C. Li, and D. Tao, "Semantic adversarial network with multi-scale pyramid attention for video classification," in *Proc. AAAI*, vol. 33, no. 01, 2019, pp. 9030–9037.
- [34] Z.-L. Ni, G.-B. Bian, G.-A. Wang, X.-H. Zhou, Z.-G. Hou, H.-B. Chen, and X.-L. Xie, "Pyramid attention aggregation network for semantic

- segmentation of surgical instruments,” in *Proc. AAAI*, vol. 34, no. 07, 2020, pp. 11 782–11 790.
- [35] C. Li, D. Du, L. Zhang, L. Wen, T. Luo, Y. Wu, and P. Zhu, “Spatial attention pyramid network for unsupervised domain adaptation,” in *Proc. ECCV*, 2020, pp. 481–497.
- [36] X. Hu, Z. Zhang, Z. Jiang, S. Chaudhuri, Z. Yang, and R. Nevatia, “Span: Spatial pyramid attention network for image manipulation localization,” in *Proc. ECCV*, 2020, pp. 312–328.
- [37] A. Anoosheh, E. Agustsson, R. Timofte, and L. Van Gool, “Combogan: Unrestrained scalability for image domain translation,” in *Proc. CVPR Workshops*, 2018, pp. 783–790.
- [38] A. Jolicoeur-Martineau, “The relativistic discriminator: a key element missing from standard gan,” *arXiv preprint arXiv:1807.00734*, 2018.
- [39] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, “Analyzing and improving the image quality of stylegan,” in *Proc. CVPR*, 2020, pp. 8110–8119.
- [40] Y. Wu and K. He, “Group normalization,” in *Proc. ECCV*, 2018, pp. 3–19.
- [41] L. I. Rudin, S. Osher, and E. Fatemi, “Nonlinear total variation based noise removal algorithms,” *Physica D: nonlinear phenomena*, vol. 60, no. 1-4, pp. 259–268, 1992.
- [42] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, “Image quality assessment: from error visibility to structural similarity,” *IEEE Trans. Image Process.*, vol. 13, no. 4, pp. 600–612, 2004.
- [43] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, “Spectral normalization for generative adversarial networks,” *arXiv preprint arXiv:1802.05957*, 2018.
- [44] K.-F. Yang, C.-Y. Li, and Y.-J. Li, “Multifeature-based surround inhibition improves contour detection in natural images,” *IEEE Trans. Image Process.*, vol. 23, no. 12, pp. 5020–5032, 2014.
- [45] M. Li, J. Liu, W. Yang, X. Sun, and Z. Guo, “Structure-revealing low-light image enhancement via robust retinex model,” *IEEE Trans. Image Process.*, vol. 27, no. 6, pp. 2828–2841, 2018.
- [46] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes challenge: A retrospective,” *International journal of computer vision*, vol. 111, no. 1, pp. 98–136, 2015.
- [47] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes (voc) challenge,” *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [48] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [49] 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 *Proc. CVPR*, 2016, pp. 3213–3223.
- [50] Z. Zheng and Y. Yang, “Unsupervised scene adaptation with memory regularization in vivo,” in *Proc. IJCAI*, 2020.
- [51] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, “Yolov4: Optimal speed and accuracy of object detection,” *arXiv preprint arXiv:2004.10934*, 2020.
- [52] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Proc. ECCV*, 2014, pp. 740–755.
- [53] M. J. Chong and D. Forsyth, “Effectively unbiased fid and inception score and where to find them,” in *Proc. CVPR*, 2020, pp. 6070–6079.