

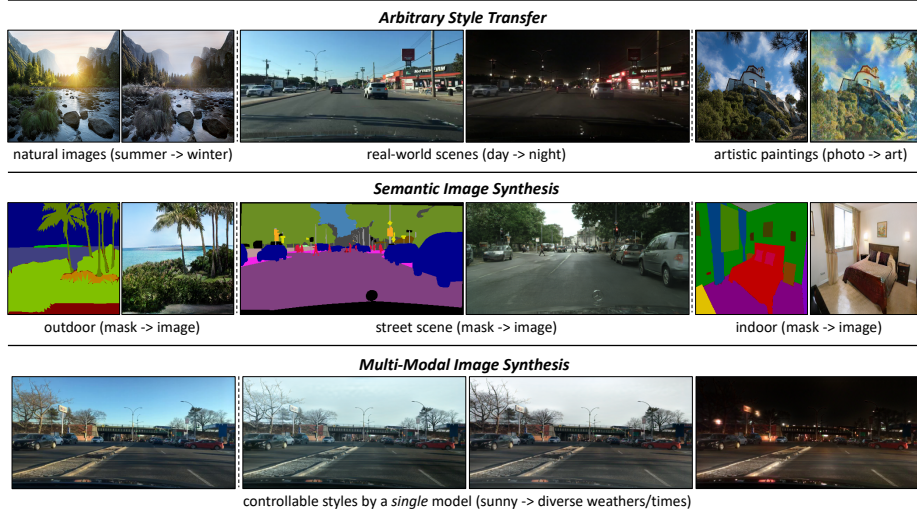
# TSIT: A Simple and Versatile Framework for Image-to-Image Translation

Liming Jiang<sup>1</sup>, Changxu Zhang<sup>2</sup>, Mingyang Huang<sup>3</sup>, Chunxiao Liu<sup>3</sup>,  
Jianping Shi<sup>3</sup>, and Chen Change Loy<sup>1</sup>✉

<sup>1</sup> Nanyang Technological University, Singapore, Singapore  
{liming002, ccloy}@ntu.edu.sg

<sup>2</sup> University of California, Berkeley, CA, USA  
zhangcx@berkeley.edu

<sup>3</sup> SenseTime Research, Beijing, China  
{huangmingyang, liuchunxiao, shijianping}@sensetime.com



**Fig. 1.** Our framework is simple and versatile for various image-to-image translation tasks. For unsupervised arbitrary style transfer, diverse scenarios (*e.g.*, natural images, real-world scenes, artistic paintings) can be handled. For supervised semantic image synthesis, our method is robust to different scenes (*e.g.*, outdoor, street scene, indoor). Multi-modal image synthesis is feasible by a *single* model with controllable styles.

**Abstract.** We introduce a simple and versatile framework for image-to-image translation. We unearth the importance of normalization layers, and provide a carefully designed two-stream generative model with newly proposed feature transformations in a coarse-to-fine fashion. This allows multi-scale semantic structure information and style representation to be

effectively captured and fused by the network, permitting our method to scale to various tasks in both unsupervised and supervised settings. No additional constraints (*e.g.*, cycle consistency) are needed, contributing to a very clean and simple method. Multi-modal image synthesis with arbitrary style control is made possible. A systematic study compares the proposed method with several state-of-the-art task-specific baselines, verifying its effectiveness in both perceptual quality and quantitative evaluations. GitHub: <https://github.com/EndlessSora/TSIT>.

## 1 Introduction

Image-to-image translation [16] aims at translating one image representation to another. Recent advances [10, 30, 21, 22, 32], especially Generative Adversarial Networks (GANs) [10], have made remarkable success in various image-to-image translation tasks. Previous studies usually present specialized solutions for a specific form of application, ranging from arbitrary style transfer [53, 44, 13, 27, 14, 24, 49] in the unsupervised setting, to semantic image synthesis [16, 4, 34, 41, 33, 28] in the supervised setting.

In this study, we are interested in devising a general and unified framework that is applicable to different image-to-image translation tasks without degradation in synthesis quality. This is *non-trivial* given the different natures of different tasks. For instance, in certain conditional image synthesis tasks (*e.g.*, arbitrary style transfer), paired data are usually not available. Under this unsupervised setting, translation task demands additional constraints on cycle consistency [53, 44, 19, 27], semantic features [39], pixel gradients [1], or pixel values [36]. In semantic image synthesis (*i.e.*, translation from segmentation labels to images), training pairs are available. This task is more data-dependent and typically needs losses to minimize per-pixel distance between the generated sample and ground truth. In addition, specialized structures [4, 41, 33, 28] are required to maintain spatial coherence and resolution. Due to the different needs, existing methods exploit their own specially designed components. It is difficult to cross-use these components or integrate them into a unified framework.

To address the aforementioned challenges, we propose a Two-Stream Image-to-image Translation (TSIT) framework, which is *versatile* for various image-to-image translation tasks (see Fig. 1). The framework is simple as it is based purely on feature transformation. Unlike previous approaches [33, 13] that only consider either semantic structure or style representation, we factorize *both* the structure and style in multi-scale *feature levels* via a symmetrical *two-stream* network. The two streams jointly influence the new image generation in a coarse-to-fine manner via a consistent feature transformation scheme. Specifically, the content spatial structure is preserved by an element-wise feature adaptive denormalization (FADE) from the content stream, while the style information is exerted by feature adaptive instance normalization (FAdaIN) from the style stream. Standard loss functions such as adversarial loss and perceptual loss are used, without additional constraints like cycle consistency. The pipeline is applicable to both unsupervised and supervised settings, easing the preparation of data.

The **contributions** of our work are summarized as follows. We propose TSIT, a simple and versatile framework, which is effective for various image-to-image translation tasks. Despite the succinct design, our network is readily adaptable to various tasks and achieves compelling results. The good performance is achieved by 1) *multi-scale* feature normalization (FADE and FAdaIN) scheme that captures *coarse-to-fine* structure and style information, and 2) a *two-stream* network design that integrates *both* content and style effectively, reducing artifacts and making multi-modal image synthesis possible (see Fig. 1). In comparison to several state-of-the-art task-specific baselines [14, 49, 4, 34, 41, 33, 28], our method achieves comparable or even better results in both perceptual quality and quantitative evaluations.

## 2 Related Work

**Image-to-image translation.** Existing methods can be classified into two categories: unsupervised and supervised. With only unpaired data, unsupervised image-to-image translation problem is inherently ill-posed. Additional constraints are needed on *e.g.*, cycle consistency [53, 44, 19, 27], semantic features [39], pixel gradients [1], or pixel values [36]. In contrast, supervised methods, such as `pix2pix` [16], are more data-dependent, requiring well-annotated paired training samples. Subsequent approaches [4, 34, 41, 33, 28] extend the supervised problem for generating high-resolution images or keeping effective semantic meaning.

Limited by learning only one-to-one mapping between two domains, some of the GAN-based methods [53, 44, 19, 27] suffer from generating images with low diversity. Recent studies explore more deeply into both multi-domain translation [6, 26] and multi-modal translation [14, 24, 48], significantly increasing generation diversity. MUNIT [14] is a representative method that disentangles domain-invariant content and domain-specific style representation, enriching the synthesized images. Multi-mapping translation is defined in a very recent work, DMIT [49], which is designed to capture multi-modal image nature in each domain.

Existing image-to-image translation methods lack the scalability to adapt to different tasks under diverse difficult settings. Different demands of unsupervised and supervised settings oblige previous methods to exploit customized modules. Cross-using these components will be suboptimal due to either degradation in quality or introduction of additional constraints. It is non-trivial to integrate them into a single framework and improve robustness. In this study, we design a two-stream network with newly proposed feature transformations inspired by [33] and [13]. Our method is succinct yet able to link various tasks.

**Arbitrary style transfer.** Style transfer is closely relevant to image-to-image translation in the unsupervised setting. Style transfer aims at retaining the content structure of an image, while manipulating its style representation adopted from other images. Classical methods [9, 17, 3, 8] gradually improve this task from optimization-based to real-time, allowing multiple style transfer during inference. Huang *et al.* introduce AdaIN [13], an effective normalization strategy for arbitrary style transfer. Several studies [45, 51, 43, 5, 23, 29, 38] improve styl-

ization via wavelet transforms [45], graph cuts [51], or iterative error-correction [38]. Besides, most collection-guided [14] style transfer methods are GAN-based [53, 44, 27, 14, 24, 49], showing impressive results.

Previous works usually consider either content or style information. In contrast, our framework succeeds in seeking a balance between content and style, and adaptively fuses them well. The proposed method achieves user-controllable multi-modal style manipulation by only a *single* model. Compared to customized style transfer methods, our approach achieves better synthesis quality in many scenarios including natural images, real-world scenes, and artistic paintings.

**Semantic image synthesis.** We define semantic image synthesis as in [33], aiming at synthesizing a photorealistic image from a semantic segmentation mask. Semantic image synthesis is a special form of supervised image-to-image translation. The domain gap of this task is large. Therefore, keeping effective semantic information to enhance fidelity without losing diversity is challenging.

**Pix2pix** [16] first adopts conditional GAN [30] in the semantic image synthesis task. Pix2pixHD [41] contains a multi-scale generator and multi-scale discriminators to generate high-resolution images. SPADE [33] takes a noise map as input, and resizes the semantic label map for modulating the activations in normalization layers by a learned affine transformation. CC-FPSE [28] employs a weight prediction network for generator. A semantics-embedding discriminator is used to enhance fine details and semantic alignments between the generated samples and the input semantic layouts. In addition to these GAN-based methods, CRN [4] applies a cascaded refinement network with regression loss as the supervision. SIMS [34] is a semi-parametric method, retrieving fragments from a memory bank and refining the canvas by a refinement network.

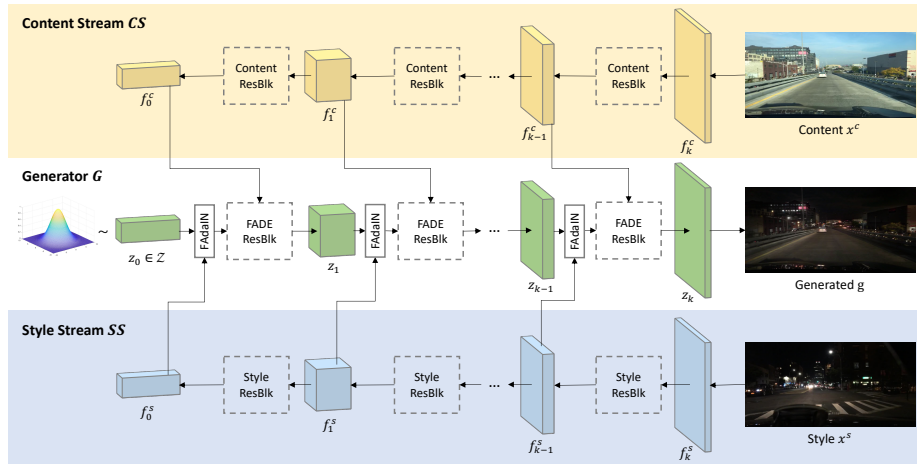
Different from prior works, we design a symmetrical two-stream framework. The network learns feature-level semantic structure information and style representation instead of directly resizing the input mask like SPADE [33]. Coarse-to-fine feature representations are learned by neural networks, adaptively keeping high fidelity without diminishing diversity.

### 3 Methodology

We consider three key requirements in formulating a robust and scalable method to link various tasks: 1) *Both* semantic structure information and style representation should be considered and fused adaptively. 2) The content and style information should be learned by networks in *feature level* instead of in image level to fit the nature of diverse semantic tasks. 3) The network structure and loss functions should be *simple* for easy training without additional constraints.

#### 3.1 Network Structure

Based on the aforementioned considerations, we design a Two-Stream Image-to-image Translation (TSIT) framework, as illustrated in Fig. 2. TSIT consists of four components: content stream, style stream, generator, and discriminators

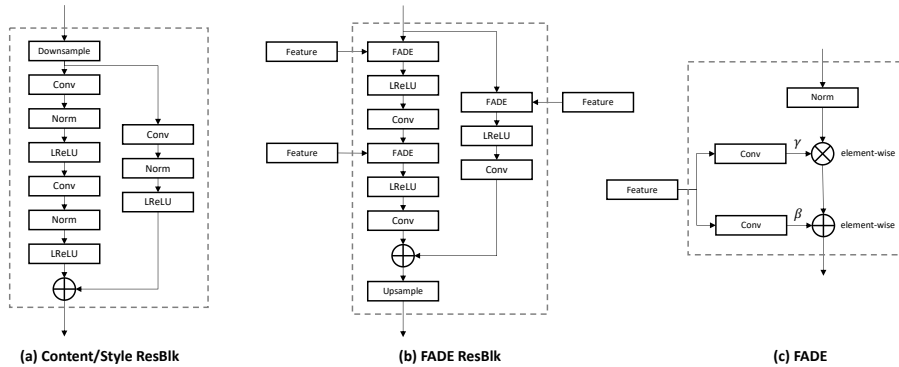


**Fig. 2.** The proposed Two-Stream Image-to-image Translation (TSIT) framework. The multi-scale patch-based discriminators are omitted. A Gaussian noise map is taken as the latent input for the generator. The feature representations of the content and style images are extracted by the corresponding streams for multi-scale feature transformations. The symmetrical networks fuse semantic structure and style representation in an end-to-end training. Submodules of our network are shown in Fig. 3.

(omitted in Fig. 2). The first three main components are fully convolutional and symmetrically designed. The details of the submodules, including content/style residual block, FADE residual block, FADE module in the FADE residual block, are as shown in Fig. 3. We will discuss them separately in this section.

**Content/style stream.** Unlike the traditional conditional GAN [30], we place the two-stream networks, *i.e.*, content stream and style stream, on each side of the generator (see Fig. 2). These two streams are symmetrical with the same network structure, aiming at extracting corresponding feature representations in different levels. We construct content/style stream based on standard residual blocks [11]. We call them content/style residual blocks. As shown in Fig 3 (a), each block has three convolutional layers, one of which is designed for the learned skip connection. The activation function is Leaky ReLU. The function of content/style stream is to extract features and feed them to the corresponding feature transformation layers in the generator. Multi-scale content/style representation in *feature levels* can be learned by content/style stream, adaptively fitting different feature transformations.

**Generator.** The generator has a completely inverse structure *w.r.t.* the content/style stream. This is intentionally designed to consistently match the level of semantic abstraction at different feature scales. A noise map is sampled from a Gaussian distribution as the latent input, and the feature maps from corresponding layers in content/style stream are taken as multi-scale feature inputs. The proposed feature transformations are implemented by a FADE residual block



**Fig. 3.** Submodules of our framework. (a) is a content/style residual block in the symmetrical content/style streams. (b) is a FADE residual block in the generator. (c) is a FADE module in the FADE residual block. It performs *element-wise* denormalization by modulating the normalized activation using a learned affine transformation defined by the modulation parameters  $\gamma$  and  $\beta$ .

(Fig. 3 (b)) and a FAdaIN module. In the FADE residual block, we use an inverse architecture *w.r.t.* the content/style residual block and replace the batch normalization [15] layer with the FADE module (Fig. 3 (c)). The FADE module performs *element-wise* denormalization by modulating the normalized activation using a learned affine transformation defined by the modulation parameters  $\gamma$  and  $\beta$ . The FAdaIN module is used to exert style information through feature adaptive instance normalization. More discussions are given in Sec. 3.2.

The entire image generation process is performed in a coarse-to-fine manner. In particular, multi-scale content/style features are injected to refine the generated image constantly from high-level latent code to low-level image representation. Semantic structure and style information are learnable and effectively fused in an end-to-end training.

**Discriminators.** We exploit the standard multi-scale patch-based discriminators (omitted in Fig. 2) in [41, 33]. Three regular discriminators with an identical architecture are included to discriminate images at different scales. Despite the same structure, patch-based training allows the discriminator operating at the coarsest scale to have the largest receptive field, capturing global information of the image. Whereas the one operating at the finest scale has the smallest receptive field, making the generator produce better details. Multi-scale patch-based discriminators further improve the robustness of our method for image-to-image translation tasks in different resolutions. Besides, the discriminators also serve as feature extractors for the generator to optimize the feature matching loss.

### 3.2 Feature Transformation

We propose a new feature transformation scheme, considering *both* semantic structure information and style representation, and fusing them adaptively. Let

$x^c$  be the content image and  $x^s$  be the style image.  $CS$ ,  $SS$ ,  $G$ ,  $D$  denote content stream, style stream, generator, and discriminators, respectively. Sampled from a Gaussian distribution,  $z_0 \in \mathbb{Z}$  is a noise map as the latent input for the generator (Fig. 2). Let  $z_i \in \{z_0, z_1, z_2, \dots, z_k\}$  be the feature map after  $i$ -th residual block in the generator, with  $k$  denoting the total number of residual blocks (*i.e.*, the upsampling times in the generator). Let  $f_i^c \in \{f_0^c, f_1^c, f_2^c, \dots, f_k^c\}$  represent the corresponding feature representations extracted by the content stream (Fig. 2),  $f_i^s \in \{f_0^s, f_1^s, f_2^s, \dots, f_k^s\}$  with the similar meaning in the style stream.

**Feature adaptive denormalization (FADE).** Our method is inspired by spatially adaptive denormalization (SPADE) [33]. Different from SPADE that resizes a semantic mask as its input, we generalize the input to multi-scale *feature representation*  $f_i^c$  of the content image  $x^c$ . In this way, we fully exploit semantic information captured by the content stream CS.

Formally, we define  $N$  as the batch size,  $L_i$  as the number of feature map channels in each layer.  $H_i$  and  $W_i$  are height and width, respectively. We first apply batch normalization [15] to normalize the generator feature map  $z_i$  in a channel-wise manner. Then, we modulate the normalized feature by using the learned parameters scale  $\gamma_i$  and bias  $\beta_i$ . The denormalized activation ( $n \in N$ ,  $l \in L_i$ ,  $h \in H_i$ ,  $w \in W_i$ ) is:

$$\gamma_i^{l,h,w} \cdot \frac{z_i^{n,l,h,w} - \mu_i^l}{\sigma_i^l} + \beta_i^{l,h,w}, \quad (1)$$

where  $\mu_i^l$  and  $\sigma_i^l$  are the mean and standard deviation, respectively, of the generator feature map  $z_i$  before the batch normalization [15] in channel  $l$ :

$$\mu_i^l = \frac{1}{NH_iW_i} \sum_{n,h,w} z_i^{n,l,h,w}, \quad (2)$$

$$\sigma_i^l = \sqrt{\frac{1}{NH_iW_i} \sum_{n,h,w} (z_i^{n,l,h,w})^2 - (\mu_i^l)^2}. \quad (3)$$

The denormalization operation is *element-wise*, and the parameters  $\gamma_i^{l,h,w}$  and  $\beta_i^{l,h,w}$  are learned by one-layer convolutions from  $f_i^c$  in the FADE module (see Fig. 3 (c)). Compared to previous conditional normalization methods [8, 13, 33], FADE experiences more perceptible influence from coarse-to-fine feature representations, thus it can better preserve semantic structure information.

**Feature adaptive instance normalization (FAdaIN).** To better fuse style representation, we introduce another feature transformation, named feature adaptive instance normalization (FAdaIN). This method is inspired by adaptive instance normalization (AdaIN) [13], with a generalization to enable the style stream  $SS$  to learn multi-scale *feature-level* style representation  $f_i^s$  of the style image  $x^s$  more effectively.

We use the same notation  $z_i$  to represent the feature map after  $i$ -th FADE residual block in the generator. FAdaIN adaptively computes the affine param-

eters from the corresponding style feature  $f_i^s$  with the same scale from  $SS$ :

$$\text{FAdaIN}(z_i, f_i^s) = \sigma(f_i^s) \left( \frac{z_i - \mu(z_i)}{\sigma(z_i)} \right) + \mu(f_i^s), \quad (4)$$

where  $\mu(z_i)$  and  $\sigma(z_i)$  are the mean and standard deviation, respectively, of  $z_i$ .

Exploiting FAdaIN, coarse-to-fine style features at different layers can be fused adaptively with the corresponding semantic structure features learned by FADE, allowing our framework to be trained end-to-end and versatile to different tasks. Furthermore, owing to the effectiveness of FAdaIN in capturing multi-scale style feature representations, multi-modal image synthesis is made possible with arbitrary style control.

### 3.3 Objective

We use standard losses in our objective function. Following [33, 28], we adopt a hinge loss term [25, 31, 50] as our adversarial loss. For the generator, we apply hinge-based adversarial loss, perceptual loss [17], and feature matching loss [41]. For the multi-scale discriminators, only hinge-based adversarial loss is used to distinguish whether the image is real or fake. The generator and discriminator are trained alternately to play a min-max game. The generator loss  $\mathcal{L}_G$  and the discriminator loss  $\mathcal{L}_D$  can be written as:

$$\mathcal{L}_G = -\mathbb{E}[D(g)] + \lambda_P \mathcal{L}_P(g, x^c) + \lambda_{FM} \mathcal{L}_{FM}(g, x^s), \quad (5)$$

$$\mathcal{L}_D = -\mathbb{E}[\min(-1 + D(x^s), 0)] - \mathbb{E}[\min(-1 - D(g), 0)], \quad (6)$$

where  $g = G(z_0, x^c, x^s)$  denotes the generated image,  $z_0$ ,  $x^c$ ,  $x^s$  denote the input noise map in latent space, the content image, and the style image, respectively.  $\mathcal{L}_P$  is the perceptual loss [17] that minimizes the difference between the feature representations extracted by VGG-19 [17] network.  $\mathcal{L}_{FM}$  is the feature matching loss [41] that matches the intermediate features at different layers of multi-scale discriminators.  $\lambda_P$  and  $\lambda_{FM}$  are the corresponding weights.

The simple objective functions make our framework stable and easy to train. Thanks to the two-stream network, the typical KL loss [21, 49, 33, 28] for multi-modal image synthesis becomes optional. Despite the simplicity, TSIT is a highly versatile tool, readily adaptable to various image-to-image translation tasks.

## 4 Settings

**Implementation details.** We use Adam [20] optimizer and set  $\beta_1 = 0$ ,  $\beta_2 = 0.9$ . Two time-scale update rule [12] is applied, where the learning rates for the generator (including two streams) and the discriminators are 0.0001 and 0.0004, respectively. We exploit Spectral Norm [31] for all layers in our network. We adopt SyncBN and IN [40] for the generator and the multi-scale discriminators, respectively. For the perceptual loss [17], we use the feature maps of



`relu1_1`, `relu2_1`, `relu3_1`, `relu4_1`, `relu5_1` layers from a pretrained VGG-19 [37] model, with the weights [1/32, 1/16, 1/8, 1/4, 1]. For the feature matching loss [41], we select features of three layers from the discriminator at each scale. All the experiments are conducted on NVIDIA Tesla V100 GPUs. Please refer to our *supplementary material* for additional implementation details.

**Applications.** The proposed framework is versatile for various image-to-image translation tasks. We consider three representative applications of conditional image synthesis: arbitrary style transfer (unsupervised), semantic image synthesis (supervised), and multi-modal image synthesis (enriching generation diversity). Please refer to our *supplementary material* for details of our application exploration.

**Datasets.** For arbitrary style transfer, we consider diverse scenarios. We use Yosemite summer  $\rightarrow$  winter dataset (natural images) provided by [53]. We classify BDD100K [47] (real-world scenes) into different times and perform day  $\rightarrow$  night translation. Besides, we use Photo  $\rightarrow$  art dataset (artistic paintings) in [53]. For semantic image synthesis, we select several challenging datasets (*i.e.*, Cityscapes [7] and ADE20K [52]). For multi-modal image synthesis, we further classify BDD100K [47] into different time and weather conditions, and perform controllable time and weather translation. The details of the datasets can be found in the *supplementary material*.

**Evaluation metrics.** Besides comparing perceptual quality, we employ the standard evaluation protocol in prior works [14, 2, 18, 33, 28] for quantitative evaluation. For arbitrary style transfer, we apply Fréchet Inception Distance (FID, evaluating similarity of distribution between the generated images and the real images, lower is better) [12] and Inception Score (IS, considering clarity and diversity, higher is better) [35]. For semantic image synthesis, we strictly follow [33, 28], adopting FID [12] and segmentation accuracy (mean Intersection-over-Union (mIoU) and pixel accuracy (accu)). The segmentation models are: DRN-D-105 [46] for Cityscapes [7], and UperNet101 [42] for ADE20K [52].

**Baselines.** We compare our method with several state-of-the-art task-specific baselines. For a fair comparison, we mainly employ GAN-based methods. In the unsupervised setting, MUNIT [14] and DMIT [49] are included, with the strong ability to capture the multi-modal nature of images while keeping quality. In the supervised setting, we compare against CRN [4], SIMS [34], pix2pixHD [41], SPADE [33], and CC-FPSE [28].

## 5 Results

**Arbitrary style transfer.** The results of *Yosemite summer  $\rightarrow$  winter season transfer* are shown in Fig. 4. Baselines [14, 49] tend to impose the color of the style image (winter) to the whole content image (summer). Besides, MUNIT sometimes introduces unnecessary artistic effects, and DMIT generates some grid-like artifacts. In comparison, our generated results are clearer and more semantics-aware spatially. The results of *BDD100K day  $\rightarrow$  night time translation* are shown in Fig. 5. Some objects (*e.g.*, road sign, car) generated by MUNIT

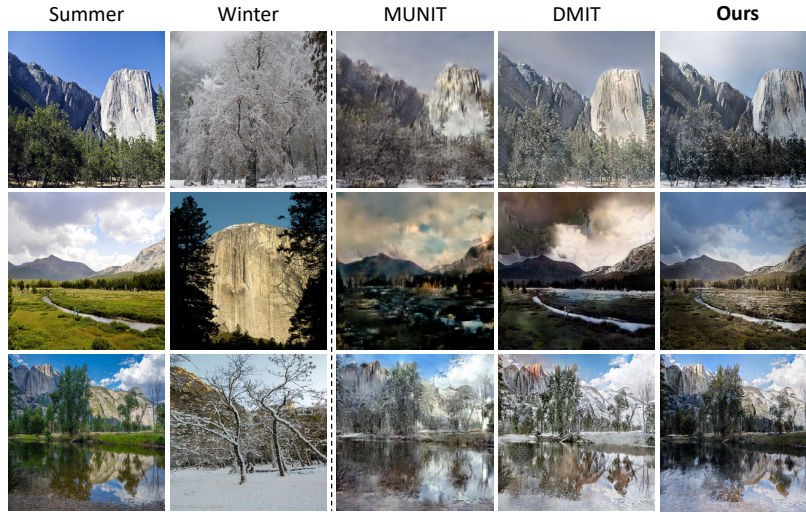


Fig. 4. Yosemite summer  $\rightarrow$  winter season transfer results compared to baselines.

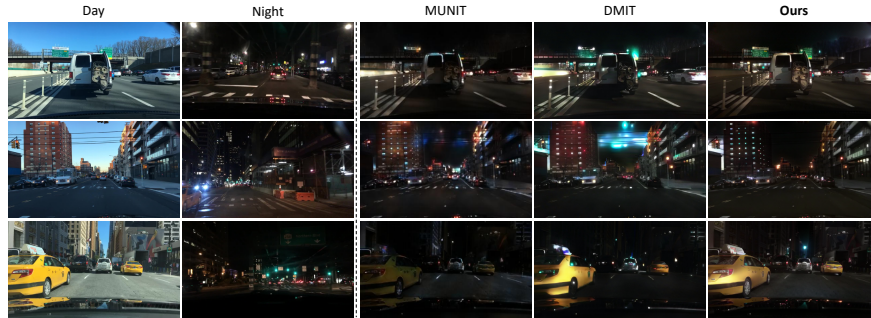


Fig. 5. BDD100K day  $\rightarrow$  night time translation results compared to baselines.

are too dark, and the whole image tends to have some unnatural colors. DMIT introduces obvious artifacts to the car or sky. In contrast, our method produces more photorealistic samples in this task. In *photo  $\rightarrow$  art style transfer*, we choose some hard cases to make a clear comparison (see Fig. 6) due to the very strong ability of all the methods in this task. Our method can transfer the styles well while effectively keeping the content structure. MUNIT tends to impose a homogeneous color to the image. Although DMIT achieves slightly better stylization than our method in certain cases (in Row 3 of Fig. 6), it also brings some grid-like distortions.

The quantitative evaluation results are shown in Table 1. Our approach achieves better performance than baselines [14, 49] in all the tasks. We also note

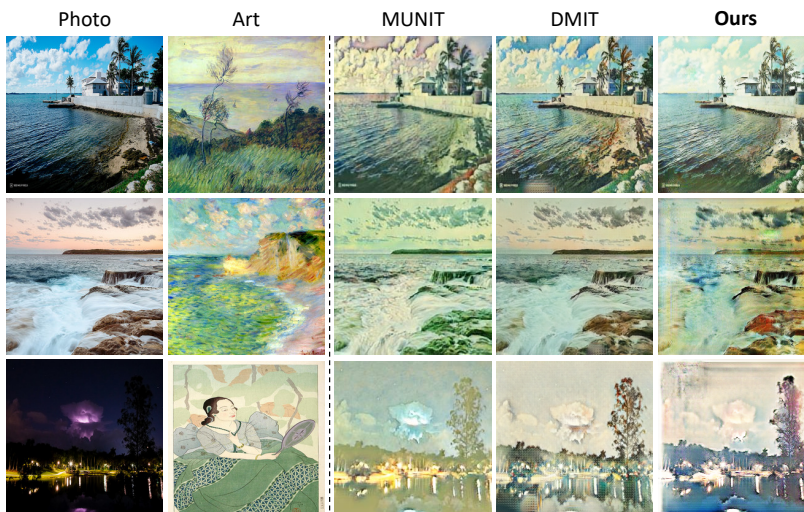


Fig. 6. Photo  $\rightarrow$  art style transfer results compared to baselines.

Table 1. The FID and IS scores of our method compared to state-of-the-art methods in arbitrary style transfer tasks. A lower FID and a higher IS indicate better performance.

Methods	summer $\rightarrow$ winter		day $\rightarrow$ night		photo $\rightarrow$ art	
	FID $\downarrow$	IS $\uparrow$	FID $\downarrow$	IS $\uparrow$	FID $\downarrow$	IS $\uparrow$
MUNIT [14]	118.225	2.537	110.011	2.185	167.314	3.961
DMIT [49]	87.969	2.884	83.898	2.156	166.933	3.871
Ours	<b>80.138</b>	<b>2.996</b>	<b>79.697</b>	<b>2.203</b>	<b>165.561</b>	<b>4.020</b>

that the gap is relatively small in photo  $\rightarrow$  art style transfer, in line with the close qualitative performance in this task (see Fig. 6).

**Semantic image synthesis.** We choose two state-of-the-art baselines, SPADE [33] and CC-FPSE [28], to show some qualitative comparison results of semantic image synthesis (Fig. 7). Our method demonstrates better perceptual quality than these task-specific baselines. In street scene (Column 1), our method generates better details on key objects (car, pedestrian). In road scene (Column 2), SPADE generates atypical colors on the roads, while CC-FPSE produces unnatural edges on the cars, hardly fitting the background (road). For outdoor natural images (Column 3), all the methods share a similar generation quality. Our method is slightly better due to less distortions on the grass. In indoor scene (Column 4 and 5), SPADE and CC-FPSE produce obvious artifacts in some cases (Column 5). In contrast, our method is more robust to diverse scenarios.

The quantitative evaluation results are shown in Table 2 (the values used for comparison are taken from [33, 28]). The proposed method achieves comparable performance with the very strong specialized methods [4, 34, 41, 33, 28] for semantic image synthesis. Note that SIMS [34] yields the best FID score but poor

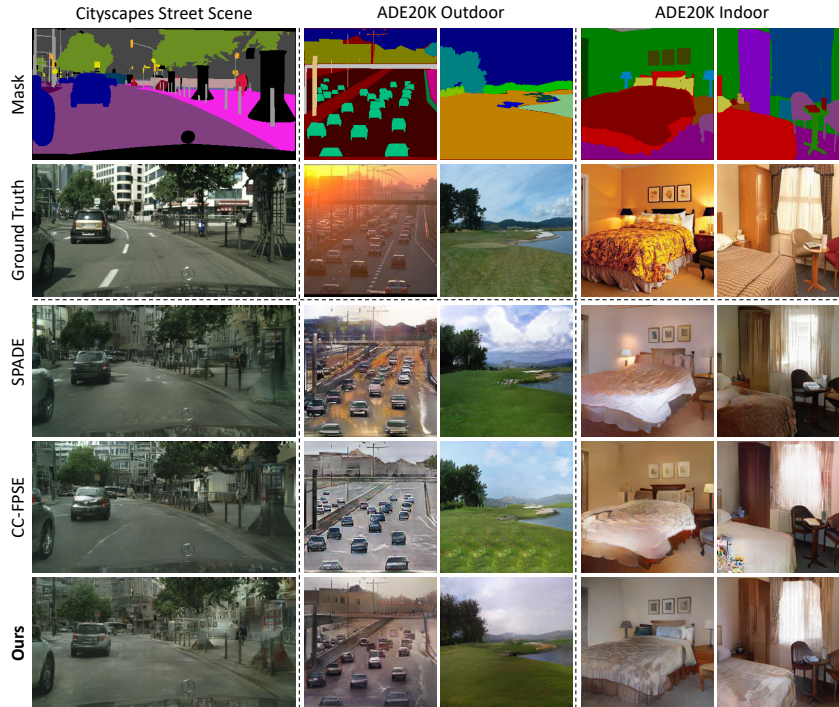


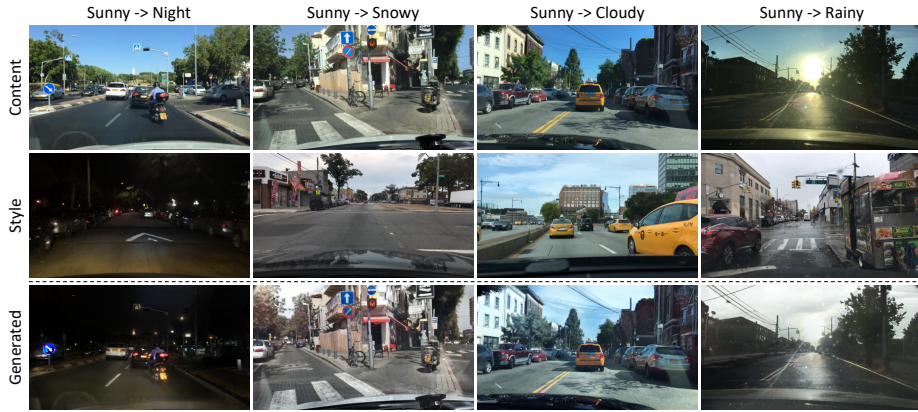
Fig. 7. Semantic image synthesis results compared to baselines.

Table 2. The mIoU, pixel accuracy (accu) and FID scores of our method compared to state-of-the-art methods in semantic image synthesis tasks. A higher mIoU, a higher pixel accuracy (accu) and a lower FID indicate better performance.

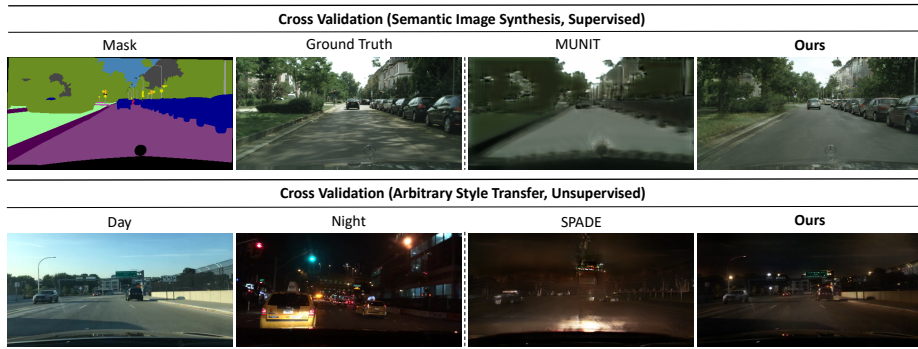
Methods	Cityscapes			ADE20K		
	mIoU $\uparrow$	accu $\uparrow$	FID $\downarrow$	mIoU $\uparrow$	accu $\uparrow$	FID $\downarrow$
CRN [4]	52.4	77.1	104.7	22.4	68.8	73.3
SIMS [34]	47.2	75.5	<b>49.7</b>	N/A	N/A	N/A
pix2pixHD [41]	58.3	81.4	95.0	20.3	69.2	81.8
SPADE [33]	62.3	81.9	71.8	38.5	79.9	33.9
CC-FPSE [28]	65.5	82.3	54.3	<b>43.7</b>	<b>82.9</b>	31.7
Ours	<b>65.9</b>	<b>94.4</b>	59.2	38.6	80.8	<b>31.6</b>

segmentation performance on Cityscapes, because it stitches image patches from a memory bank of training set while not keeping the exactly consistent position in the synthesized image. Our approach achieves state-of-the-art segmentation performance on Cityscapes and the best FID score on ADE20K, suggesting its robustness to fit the nature of different image-to-image translation tasks.

**Multi-modal image synthesis.** We perform multi-modal image synthesis for time and weather image-to-image translation (see Fig. 8) on BDD100K [47].



**Fig. 8. BDD100K multi-modal image synthesis** for different time and weather translation results by a *single* model.



**Fig. 9. Cross validation** of ineffectiveness of task-specific methods in inverse settings.

Training only a *single* model, we translate the images of weather *sunny* to different times and weathers (*i.e.*, *night*, *snowy*, *cloudy*, *rainy*). Our method effectively adapts to different style control and keeps photorealistic generation quality. Although the weather *snowy* is not very obvious in BDD100K [47], our approach successfully introduces some snow-like effects on trees and grounds (Column 2). **Cross validation.** We also conduct experiments to evaluate the performance of existing specialized methods in inverse settings (*i.e.*, using unsupervised methods to do semantic image synthesis / using supervised methods to perform arbitrary style transfer). We selected two representative methods, MUNIT [14] and SPADE [33]. Without modifying the architecture, we tuned the loss weights and tried to get the best generation results. To ensure a fair comparison, we also tried to compute perceptual loss with the content (day) image for SPADE to match the setting of TSIT. Representative results of cross validation are shown in Fig. 9. The proposed method shows much better results than baseline methods. MUNIT



**Fig. 10. Ablation studies** of key modules (*i.e.*, content stream (CS), style stream(SS)) and feature transformations in multi-modal image synthesis task.

fails to adapt to semantic image synthesis. SPADE loses details of key objects and introduces very strong artifacts despite translating the color correctly.

**Ablation studies.** We present ablation studies of key modules (*i.e.*, content stream (CS), style stream(SS)) and the proposed feature transformations (see Fig. 10. More ablation study results can be found in the *supplementary material*). We perform multi-modal image synthesis to show the effectiveness of different components. Our full model generates high-quality results (Column 3). When we directly inject the resized content image without CS, the semantic structure information becomes weak, leading to several artifacts in the sky (Column 4). Without SS, the model cannot perform multi-modal image synthesis at all (Column 5). The style representation is dominated by the night style. When we concatenate the feature maps of CS with the ones of the generator instead of using FADE, the concatenation introduces too much content information, leading to several failure cases (*e.g.*, *sunny*  $\rightarrow$  *night* in Column 6). If we discard FAdaIN by concatenating the feature maps of SS with the ones of the generator, the style becomes too strong, causing serious style regionalization problem (Column 7).

## 6 Conclusion

We propose TSIT, a simple and versatile framework for image-to-image translation. The proposed symmetrical two-stream network allows the image generation to be effectively conditioned on the multi-scale feature-level semantic structure information and style representation via feature transformations. A systematic study verifies the effectiveness of our method in diverse tasks compared to state-of-the-art task-specific baselines. We believe that designing a unified and versatile framework for more tasks is an important direction in the image generation area. Incorporating unconditional image synthesis tasks and introducing more variability into the two streams/latent space can be interesting future works.

**Acknowledgements.** This work is supported by the SenseTime-NTU Collaboration Project, Singapore MOE AcRF Tier 1 (2018-T1-002-056), and NTU NAP.

## References

1. Bousmalis, K., Silberman, N., Dohan, D., Erhan, D., Krishnan, D.: Unsupervised pixel-level domain adaptation with generative adversarial networks. In: CVPR (2017)
2. Brock, A., Donahue, J., Simonyan, K.: Large scale GAN training for high fidelity natural image synthesis. In: ICLR (2018)
3. Chen, D., Yuan, L., Liao, J., Yu, N., Hua, G.: StyleBank: An explicit representation for neural image style transfer. In: CVPR (2017)
4. Chen, Q., Koltun, V.: Photographic image synthesis with cascaded refinement networks. In: ICCV (2017)
5. Chiu, T.Y.: Understanding generalized whitening and coloring transform for universal style transfer. In: ICCV (2019)
6. Choi, Y., Choi, M., Kim, M., Ha, J.W., Kim, S., Choo, J.: StarGAN: Unified generative adversarial networks for multi-domain image-to-image translation. In: CVPR (2018)
7. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The Cityscapes dataset for semantic urban scene understanding. In: CVPR (2016)
8. Dumoulin, V., Shlens, J., Kudlur, M.: A learned representation for artistic style. arXiv preprint [arXiv:1610.07629](https://arxiv.org/abs/1610.07629) (2016)
9. Gatys, L.A., Ecker, A.S., Bethge, M.: A neural algorithm of artistic style. arXiv preprint [arXiv:1508.06576](https://arxiv.org/abs/1508.06576) (2015)
10. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: NeurIPS (2014)
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
12. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: GANs trained by a two time-scale update rule converge to a local nash equilibrium. In: NeurIPS (2017)
13. Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: ICCV (2017)
14. Huang, X., Liu, M.Y., Belongie, S., Kautz, J.: Multimodal unsupervised image-to-image translation. In: ECCV (2018)
15. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: ICML (2015)
16. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: CVPR (2017)
17. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: ECCV (2016)
18. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: CVPR (2019)
19. Kim, T., Cha, M., Kim, H., Lee, J.K., Kim, J.: Learning to discover cross-domain relations with generative adversarial networks. In: ICML (2017)
20. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint [arXiv:1412.6980](https://arxiv.org/abs/1412.6980) (2014)
21. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint [arXiv:1312.6114](https://arxiv.org/abs/1312.6114) (2013)
22. Kingma, D.P., Mohamed, S., Rezende, D.J., Welling, M.: Semi-supervised learning with deep generative models. In: NeurIPS (2014)

23. Kotovenko, D., Sanakoyeu, A., Lang, S., Ommer, B.: Content and style disentanglement for artistic style transfer. In: ICCV (2019)
24. Lee, H.Y., Tseng, H.Y., Huang, J.B., Singh, M., Yang, M.H.: Diverse image-to-image translation via disentangled representations. In: ECCV (2018)
25. Lim, J.H., Ye, J.C.: Geometric GAN. arXiv preprint **arXiv:1705.02894** (2017)
26. Liu, A.H., Liu, Y.C., Yeh, Y.Y., Wang, Y.C.F.: A unified feature disentangler for multi-domain image translation and manipulation. In: NeurIPS (2018)
27. Liu, M.Y., Breuel, T., Kautz, J.: Unsupervised image-to-image translation networks. In: NeurIPS (2017)
28. Liu, X., Yin, G., Shao, J., Wang, X., et al.: Learning to predict layout-to-image conditional convolutions for semantic image synthesis. In: NeurIPS (2019)
29. Lu, M., Zhao, H., Yao, A., Chen, Y., Xu, F., Zhang, L.: A closed-form solution to universal style transfer. In: ICCV (2019)
30. Mirza, M., Osindero, S.: Conditional generative adversarial nets. arXiv preprint **arXiv:1411.1784** (2014)
31. Miyato, T., Kataoka, T., Koyama, M., Yoshida, Y.: Spectral normalization for generative adversarial networks. arXiv preprint **arXiv:1802.05957** (2018)
32. Van den Oord, A., Kalchbrenner, N., Espeholt, L., Vinyals, O., Graves, A., et al.: Conditional image generation with PixelCNN decoders. In: NeurIPS (2016)
33. Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. In: CVPR (2019)
34. Qi, X., Chen, Q., Jia, J., Koltun, V.: Semi-parametric image synthesis. In: CVPR (2018)
35. Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X.: Improved techniques for training GANs. In: NeurIPS (2016)
36. Shrivastava, A., Pfister, T., Tuzel, O., Susskind, J., Wang, W., Webb, R.: Learning from simulated and unsupervised images through adversarial training. In: CVPR (2017)
37. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint **arXiv:1409.1556** (2014)
38. Song, C., Wu, Z., Zhou, Y., Gong, M., Huang, H.: ETNet: Error transition network for arbitrary style transfer. In: NeurIPS (2019)
39. Taigman, Y., Polyak, A., Wolf, L.: Unsupervised cross-domain image generation. arXiv preprint **arXiv:1611.02200** (2016)
40. Ulyanov, D., Vedaldi, A., Lempitsky, V.: Instance normalization: The missing ingredient for fast stylization. arXiv preprint **arXiv:1607.08022** (2016)
41. Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B.: High-resolution image synthesis and semantic manipulation with conditional GANs. In: CVPR (2018)
42. Xiao, T., Liu, Y., Zhou, B., Jiang, Y., Sun, J.: Unified perceptual parsing for scene understanding. In: ECCV (2018)
43. Yang, S., Wang, Z., Wang, Z., Xu, N., Liu, J., Guo, Z.: Controllable artistic text style transfer via shape-matching GAN. In: ICCV (2019)
44. Yi, Z., Zhang, H., Tan, P., Gong, M.: DualGAN: Unsupervised dual learning for image-to-image translation. In: ICCV (2017)
45. Yoo, J., Uh, Y., Chun, S., Kang, B., Ha, J.W.: Photorealistic style transfer via wavelet transforms. In: ICCV (2019)
46. Yu, F., Koltun, V., Funkhouser, T.: Dilated residual networks. In: CVPR (2017)
47. Yu, F., Xian, W., Chen, Y., Liu, F., Liao, M., Madhavan, V., Darrell, T.: BDD100K: A diverse driving video database with scalable annotation tooling. arXiv preprint **arXiv:1805.04687** (2018)



48. Yu, X., Cai, X., Ying, Z., Li, T., Li, G.: SingleGAN: Image-to-image translation by a single-generator network using multiple generative adversarial learning. In: ACCV (2018)
49. Yu, X., Chen, Y., Liu, S., Li, T., Li, G.: Multi-mapping image-to-image translation via learning disentanglement. In: NeurIPS (2019)
50. Zhang, H., Goodfellow, I., Metaxas, D., Odena, A.: Self-attention generative adversarial networks. arXiv preprint **arXiv:1805.08318** (2018)
51. Zhang, Y., Fang, C., Wang, Y., Wang, Z., Lin, Z., Fu, Y., Yang, J.: Multimodal style transfer via graph cuts. In: ICCV (2019)
52. Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralba, A.: Scene parsing through ADE20K dataset. In: CVPR (2017)
53. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: ICCV (2017)