

# Diverse Semantic Image Synthesis via Probability Distribution Modeling

Zhentaο Tan<sup>1</sup>, Menglei Chai<sup>2</sup>, Dongdong Chen<sup>3</sup>, Jing Liao<sup>4</sup>,  
Qi Chu<sup>1</sup>, Bin Liu<sup>1</sup>, Gang Hua<sup>5</sup>, Nenghai Yu<sup>1</sup>

<sup>1</sup>University of Science and Technology of China <sup>2</sup>Snap Inc. <sup>3</sup>Microsoft Cloud AI

<sup>4</sup>City University of Hong Kong <sup>5</sup>Wormpex AI Research LLC

{tzt@mail., qchu@, flowice@, ynh@}ustc.edu.cn

mchai@snap.com, {cddlyf, ganghua}@gmail.com, jingliao@cityu.edu.hk

## Abstract

*Semantic image synthesis, translating semantic layouts to photo-realistic images, is a one-to-many mapping problem. Though impressive progress has been recently made, diverse semantic synthesis that can efficiently produce semantic-level multimodal results, still remains a challenge. In this paper, we propose a novel diverse semantic image synthesis framework from the perspective of semantic class distributions, which naturally supports diverse generation at semantic or even instance level. We achieve this by modeling class-level conditional modulation parameters as continuous probability distributions instead of discrete values, and sampling per-instance modulation parameters through instance-adaptive stochastic sampling that is consistent across the network. Moreover, we propose prior noise remapping, through linear perturbation parameters encoded from paired references, to facilitate supervised training and exemplar-based instance style control at test time. Extensive experiments on multiple datasets show that our method can achieve superior diversity and comparable quality compared to state-of-the-art methods. Code will be available at <https://github.com/tzt101/INADE.git>*

## 1. Introduction

Image synthesis has recently seen impressive progress, particularly with the help of generative adversarial networks (GANs) [8]. Besides stochastic approaches that generate high-quality images from random latent variables [18, 19], conditional image synthesis is attracting equal or even more attention due to the practical advantages of its controllability. The conditional input, to guide the synthesis, can be of various forms, including RGB images, edge/gradient maps, semantic labels, etc. In this work, semantic image synthesis is one particular task that aims to generate a photo-realistic

image from a semantic label mask. In particular, we further explore its diversity and controllability without loss of generation quality. Some samples are shown in Figure 1.

Previous works [15, 42] propose solutions within the general image-to-image translation framework, which directly feeds the semantic mask into the encoder-decoder network. For higher quality, some recent methods [30, 51, 36] adopt spatially-varying conditional normalization to avoid the loss of semantic information due to conventional normalization layers [40]. Although proven successful in synthesizing certain types of content, these methods lack controllability over the generation diversity, which is particularly important for such a one-to-many problem. Some methods [50, 44] attempt to yield multimodal results by incorporating variational auto-encoder (VAE) or introducing noises. However, these methods only support global image-level diversity. To obtain finer-grained controllability, a recent work [52] proposes to use group convolution for different semantics to achieve semantic-level diversity. However, it is computationally expensive and difficult to be extended to support diversity at the instance level.

In this paper, we attempt to achieve controllable diversity in semantic image synthesis from the perspective of semantic probability distributions. The intuition is to treat each semantic class as one distribution, so that each instance of this class could be drawn from this distribution as a discrete sample. Following this idea, we propose a novel semantic image synthesis framework, which is naturally capable of producing diverse results at semantic or even instance level.

Specifically, our method contains three key ingredients. Firstly, we propose *variational modulation models* (§ 3.2) that extend discrete modulation parameters to class-wise continuous probability distributions, which embed diverse styles of each semantic category in a class-adaptive manner. Secondly, based on the variational models built per normalization layer, we further develop an *instance-adaptive sampling* method (§ 3.3) that achieves instance-level diversity by stochastically sampling modulation parameters from the

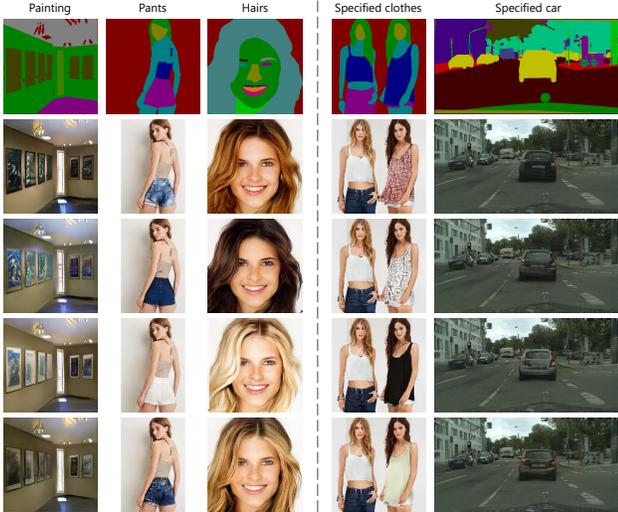


Figure 1. Semantic-level (left three columns) and Instance-level (right two columns) multimodal images generated by the proposed method. Text of each column indicates which semantic class or instance will be changed in the following results.

variational models. We harmonize the sampling across the network via consistent randomness and a learnable transformation function for each normalization layer. Finally, to more efficiently embed the instance diversity to the modulation models, we propose *prior noise remapping* (§ 3.4) that transforms the noise samples with perturbation parameters encoded from arbitrary references. We adopt this step to facilitate supervised training and enable test-time reference-based style guidance. Inspired by [30, 38, 37], the proposed method is called INADE (INstance-Adaptive DENormalization).

To evaluate the proposed method, we conduct extensive experiments on multiple datasets, including Cityscapes [3], ADE20K [48], CelebAMask-HQ [23, 17, 29], and DeepFashion [28]. Both quantitative and qualitative results show that our method significantly outperforms state-of-the-art methods by achieving much better instance-level diversity while keeping comparable generation quality.

## 2. Related Work

### 2.1. Conditional Image Synthesis

Conditional image synthesis aims at generating photo-realistic images conditioned on different types of input. We are interested in a special form of it, called semantic image synthesis, which takes segmentation layouts as input. Many impressive works have been proposed for this task. The most representative work, Pix2Pix [15] adopts an encoder-decoder generator for unified image-to-image translation. Pix2pixHD [42] improves Pix2Pix by proposing coarse-to-fine generator and discriminators. Subsequent meth-

ods [32, 27, 39, 46, 37, 51] further explore how to synthesize high quality images from semantic masks and achieve significant improvements. Besides using class-level semantic masks, some works also consider instance-level information for image synthesis, since the semantic mask itself does not provide sufficient information to synthesize instances especially in complex environments with multiple of them interacting with each other. Some works [42, 30, 37] extract boundary information from the instance map and concatenate it with the semantic mask. While recent work [6] proposes to use the instance map to guide convolution and upsampling layers for better exploiting both semantic and instance information. Different from these methods, we are interested in taking full advantage of information from instance maps to achieve instance-level diversity control.

### 2.2. Diversity in Image Synthesis

Diversity is a core target for image synthesis, which aims to generate multiple possible outputs from a single input image. Early conditional image synthesis networks either trained with paired data, like Pix2Pix [15] and Pix2pixHD [42], or with unpaired data, like CycleGAN [49], DiscoGAN [20] and UNIT [26], are single-modal. They produce one single output conditioned solely on an input image. Later, some multimodal unpaired image synthesis networks [13, 24, 1] are proposed. However, constrained by the reconstruction loss, the semantic image synthesis task trained with paired data is more difficult to support diversity. Simply concatenating a random noise vector to the input segmentation mask is usually not effective, because the generator often ignores the additional noise vectors and mode collapse may occur easily. To tackle this problem, BicycleGAN [50] enforces the bijection mapping between the noise vector and target domain. DSCGAN [44] proposes a simple regularization which can be easily integrated into most conditional GAN objectives. More recently, a variational autoencoder architecture is used to handle multimodal synthesis by [30, 37, 27]. However, these multimodal image synthesis networks only support diversity at the global level. To further control the diversity at the semantic level, the method proposed by [9] builds several auto-encoders for each face component to extract different component representations. GroupDNet [52] unifies the generation process in only one model, but still requires high computing resources, and the use of group convolution layer makes it difficult to extend to the instance level. In contrast, we propose a novel instance-aware conditional normalization framework that allows diverse instance-level generation with less overhead.

## 3. Method

We are interested in the task of semantic image synthesis, which is defined as to map a semantic mask  $\mathbf{m} \in \mathbb{L}_m^{H \times W}$

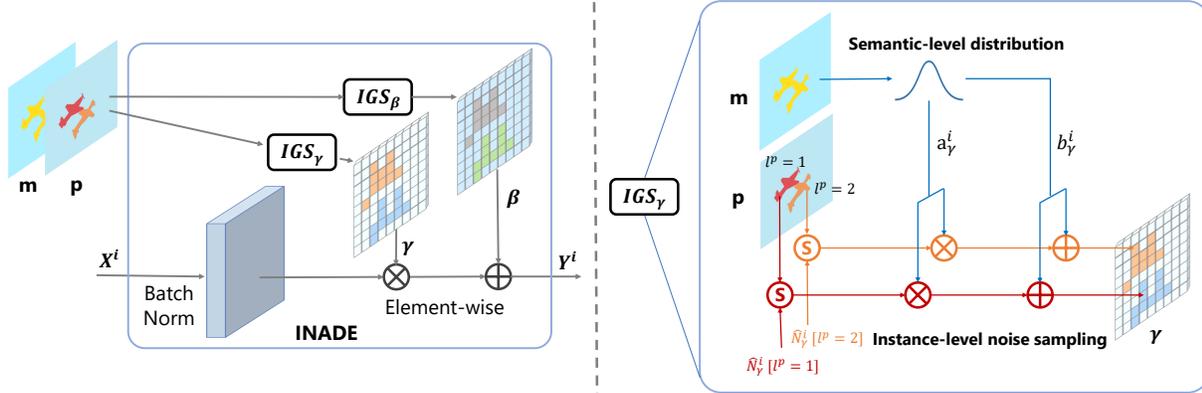


Figure 2. The illustration diagram of the proposed INstance-Adaptive DENormalization (INADE). It combines semantic-level distribution modeling and instance-level noise sampling. IGS denotes the Instance Guided Sampling which is similar to the guided sampling in [37].

to a photo-realistic image  $\mathbf{o} \in \mathbb{R}^{3 \times H \times W}$ . Here,  $\mathbf{m}$  is a class-level label map with each pixel representing an integer index to a pre-defined set of semantic categories  $\mathbb{L}_m = \{1, 2, \dots, L^m\}$ . Each pair of input  $\mathbf{m}$  and output  $\mathbf{o}$  is spatially-aligned and of the same dimension  $H \times W$ , so that the synthesized content in  $\mathbf{o}$  should comply with the corresponding semantic labels in  $\mathbf{m}$ .

In addition to this basic formulation [42, 30, 37], instance-aware semantic image synthesis [6] adopts the instance map  $\mathbf{p} \in \mathbb{L}_p^{H \times W}$  as an extra input, which differentiates different object instances sharing a same semantic label by denoting each individual instance with a unique index from the instance label set  $\mathbb{L}_p = \{1, 2, \dots, L^p\}$  in the image. By enforcing an identical semantic label within each instance, pixels belonging to a same instance label  $l^p$  in  $\mathbf{p}$  should always have a same semantic label  $l^m$  in  $\mathbf{m}$ . We represent instance to semantic label mapping as a function  $l^m = \mathcal{G}(l^p)$ .

Overall, image synthesis with instance information can be basically formulated as a function  $\mathcal{T}(\mathbf{m}, \mathbf{p}) : (\mathbb{L}_m, \mathbb{L}_p) \rightarrow \mathbb{R}^3$ . And feed-forward image translation neural networks, trained in a supervised manner, can be used to model this function. In the following, we introduce the proposed method with both the inputs of semantic and instance maps. When there is no instance label,  $\mathbf{p}$  degenerates into  $\mathbf{m}$ . And the diversity of the synthesized images changes from the instance level to the semantic level.

### 3.1. Conditional Normalization

Our solution to semantic image synthesis is based on a novel instance-level conditional normalization method. Before introducing our method, here we give a brief overview of the general framework of conditional normalization first.

Let  $\mathbf{X}^i \in \mathbb{R}^{C^i \times H^i \times W^i}$  be the activation tensor to the  $i$ -th normalization layer, where  $C^i, H^i, W^i$  are the channel depth, height, and width, respectively. In the channel-wise normalization framework similar to [14], we can generally

formulate the normalization operations as two steps: In the normalization step,  $\mathbf{X}^i$  is normalized to  $\hat{\mathbf{X}}^i$  by channel-wise mean and standard deviation  $\{\mu^i, \sigma^i\} \in \mathbb{R}^{C^i}$  in the mini-batch containing  $\mathbf{X}^i$ . Then, the modulation step scales and translates  $\hat{\mathbf{X}}^i$  with learned modulation parameters  $\{\gamma^i, \beta^i\} \in \mathbb{R}^{C^i \times H^i \times W^i}$ , which are not necessarily channel-wise constant. Let  $\mathbf{Y}^i$  be the output, for each element ( $k \in C^i, x \in H^i, y \in W^i$ ) in the tensor, we have:

$$\begin{aligned} \hat{\mathbf{X}}_{k,x,y}^i &= (\mathbf{X}_{k,x,y}^i - \mu_k^i) / \sigma_k^i, \\ \mathbf{Y}_{k,x,y}^i &= \gamma_{k,x,y}^i \hat{\mathbf{X}}_{k,x,y}^i + \beta_{k,x,y}^i. \end{aligned} \quad (1)$$

For conditional normalization [5, 12], the modulation parameters  $\gamma^i$  and  $\beta^i$  are learned with extra conditions. Specifically, for semantic image synthesis, the modulation is usually conditioned on the semantic mask  $\mathbf{m}$  [30, 37].

### 3.2. Variational Modulation Model

Conditional normalization (§ 3.1), either spatially-adaptive [30] or class-adaptive [37], has been proven helpful for semantic image synthesis. The semantic-conditioned modulation is able to largely prevent the “wash-away” effect of semantic information caused by repetitive normalizations. However, challenges still exist to achieve promising generation results with semantic-level or even instance-level diversity, given that normalization is solely conditioned on the semantic map and only global randomness is used to diversify the image styles [30]. Semantic-level diversity is realized by [52] through group convolution, but using this convolution cuts off the possibility of its extension to instance-level diversity through instance map. Recent efforts on instance-aware synthesis [42, 6] are majorly focused on better object boundaries, but not the diversity and realism of each individual instance. Due to the lack of proper instance conditioning, existing methods tend to converge instances with the same semantic label into a similar style, which significantly harms the diversity of generation.

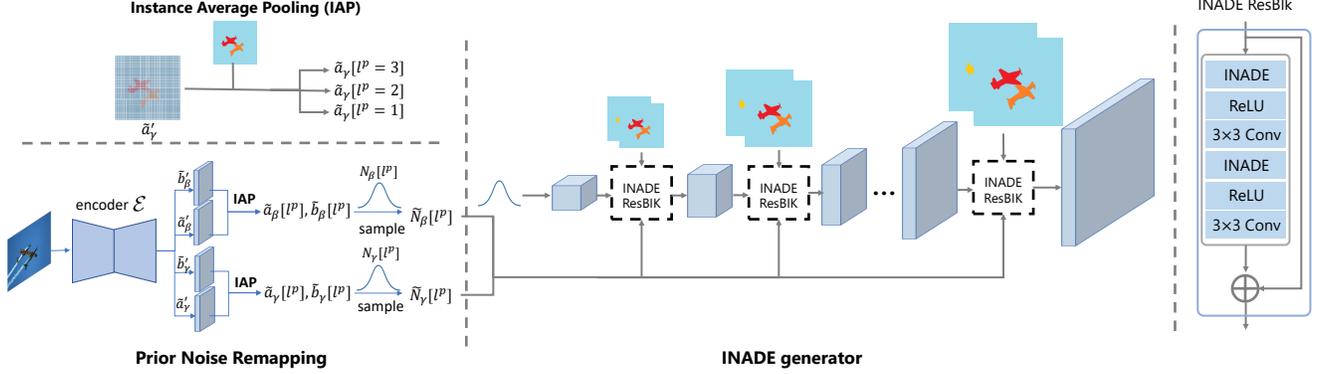


Figure 3. The overall framework of the proposed INADE generator, which consists of a remapping encoder  $\mathcal{E}$  and INADE generator.  $\mathcal{E}$  is used to transform the noise sample based on arbitrary references (§ 3.4), while the generator consists of several INADE ResBlks.

The key to instance-level diversity is a proper combination of uniform semantic-level distributions that deterministically decide the general features of a particular semantic label, and instance-level randomness that introduces allowed diversity covered by the semantic distribution models. Therefore, we model the modulation parameters as parametric probability distributions for each semantic label  $l^m \in \mathbb{L}^m$ , instead of discrete values. With such a, namely, variational modulation model, given an instance  $l^p \in \mathbb{L}^p$ , instance-level diversity is achievable via sampling modulation parameters from the probability distributions of the corresponding semantic label  $\mathcal{G}(l^p)$ . For the sake of simplicity and efficiency, following [37], we make the modulation parameters spatially-invariant and only depend on the local instance labels.

Specifically, for each semantic category  $l^m$ , its channel-wise modulation parameters are modeled as learnable probability distributions, which are built for each normalization layer in the network respectively. Formally, for the  $i$ -th layer with channel depth  $C^i$ , we have  $\{\mathbf{a}_\gamma^i, \mathbf{b}_\gamma^i, \mathbf{a}_\beta^i, \mathbf{b}_\beta^i\} \in \mathbb{R}^{L^m \times C^i}$  as the distribution transformation parameters of  $\gamma$  and  $\beta$ , respectively. All of them are treated as learnable parameters that are jointly trained with the network. Given stochastic noise matrices  $\{\mathbf{N}_\gamma^i, \mathbf{N}_\beta^i\} \in \mathbb{R}^{L^p \times C^i}$  from the same distribution for sampling, the corresponding modulation parameters of one instance label  $l^p$  in  $\mathbf{p}$  are:

$$\begin{aligned} \gamma^i[l^p] &= \mathbf{a}_\gamma^i[\mathcal{G}(l^p)] \otimes \mathbf{N}_\gamma^i[l^p] + \mathbf{b}_\gamma^i[\mathcal{G}(l^p)], \\ \beta^i[l^p] &= \mathbf{a}_\beta^i[\mathcal{G}(l^p)] \otimes \mathbf{N}_\beta^i[l^p] + \mathbf{b}_\beta^i[\mathcal{G}(l^p)], \end{aligned} \quad (2)$$

where  $\otimes$  represents element-wise multiplication, and  $[\cdot]$  accesses the vector from a matrix in the row-major order.

### 3.3. Instance-Adaptive Modulation Sampling

Our multimodal synthesis method follows the basic form of conditional modulation (§ 3.1), but further extends the conditional inputs to include not just the segmentation mask

$\mathbf{m}$ , but also the instance map  $\mathbf{p}$  and random noises to initiate sampling, as shown in Figure 2. Utilizing our variational modulation models (§ 3.2), we are able to generate diverse modulation parameters obeying the same set of probability distributions. However, considering that the generation network contains multiple conditional normalization layers, a unified sampling solution is still essential to harmonize all these layers. A straight-forward approach, independent stochastic sampling for each normalization layer, could potentially introduce inconsistency and cause the diversity to be severely neutralized. Therefore, in this paper, we propose an instance-adaptive modulation sampling method that achieves consistent instance sampling across multiple normalization layers with unequal channel depths.

To initialize, for each layer  $i$ , we resize and convert each input pair of semantic mask  $\mathbf{m}$  and instance map  $\mathbf{p}$  into the one-hot format as  $\mathbf{M}^i \in \mathbb{B}^{L^m \times H^i \times W^i}$  and  $\mathbf{P}^i \in \mathbb{B}^{L^p \times H^i \times W^i}$ , respectively, which will then be used as the conditional inputs to that layer, as shown in Figure 8. Here,  $\mathbb{B}$  represents the Boolean domain, and  $L^m, L^p$  are the aforementioned total numbers of semantic/instance labels.

For the sake of simplicity, since scale  $\gamma^i$  and shift  $\beta^i$  are generated similarly and independently, without loss of generality, here we take scale  $\gamma^i$  as the example. The sampling contains the following steps.

First of all, random samples  $\mathbf{N}_\gamma \in \mathbb{R}^{L^p \times C^0}$  are independently sampled from the standard normal distribution:  $\mathbf{N}_\gamma \sim \mathcal{N}(0, 1)$ . We use the same set of random noise samples  $\mathbf{N}_\gamma$  for all instances of normalization layers in the network, which helps enforce consistent instance styles throughout the network. Here  $C^0$  is a hyper-parameter that defines the number of the initial sampling channels.

To sample the modulation parameters for each normalization layer  $i$ , we translate the initial samples  $\mathbf{N}_\gamma$  to  $\hat{\mathbf{N}}_\gamma^i$  with a learnable linear transformation mapping  $\mathcal{F}_\gamma^i: \mathbb{R}^{L^p \times C^0} \rightarrow \mathbb{R}^{L^p \times C^i}$ :

$$\hat{\mathbf{N}}_\gamma^i = \mathcal{F}_\gamma^i(\mathbf{N}_\gamma), \quad (3)$$

where  $C^i$  is exactly the channel depth of  $i$ -th activations  $\mathbf{X}^i$ , so that the output  $\hat{N}_\gamma^i \in \mathbb{R}^{L^p \times C^i}$  assigns a transformed sample for each instance per each channel, in a spatially-invariant manner. Thus, the same source of randomness helps achieve style consistency, while the learnable transformations enforce compatible target dimensions and preserve certain ability to adapt the samples for each layer.

Finally, given the distribution transformation parameters  $\mathbf{a}_\gamma^i, \mathbf{b}_\gamma^i$  and transformed noise samples  $\hat{N}_\gamma^i$ , the scale parameters  $\gamma^i$  are calculated with Equation 2. And similarly for the shift parameters  $\beta^i$ .

### 3.4. Prior Noise Remapping

While our variational modulation models (§ 3.2) help achieve instance-level diversity, the noises, sampled regardless of instance styles (§ 3.3), can potentially introduce ambiguities during the supervised training (especially for the popular perceptual and feature matching losses), since similar noise samples can possibly correspond to instances of distinct styles. This will affect the effective diversity of generated instances and prohibit the possibility to control the instance styles with certain references.

In light of this, we propose a prior noise remapping step, during which a set of linear perturbation parameters are encoded from given references, to remap the noise samples while preserving the original distribution, in order to provide guidance to embed more meaningful instance diversity in the modulation models. To achieve this, we adopt a noise remapping encoder  $\mathcal{E}(\mathbf{r})$  that translates the reference image  $\mathbf{r} \in \mathbb{R}^{H \times W}$  into four perturbation maps  $\{\tilde{\mathbf{a}}'_\gamma, \tilde{\mathbf{b}}'_\gamma, \tilde{\mathbf{a}}'_\beta, \tilde{\mathbf{b}}'_\beta\} \in \mathbb{R}^{H \times W}$ , which are per-pixel linear transformation parameters, including scale  $\tilde{\mathbf{a}}$  and shift  $\tilde{\mathbf{b}}$ , for both  $N_\gamma$  and  $N_\beta$ . A instance aware partial convolution [10, 25] is used to avoid information contamination between different instances. Based on these dense perturbation maps, we apply an instance average pooling layer to each of these maps to get the instance-wise perturbation parameters  $\{\tilde{\mathbf{a}}_\gamma, \tilde{\mathbf{b}}_\gamma, \tilde{\mathbf{a}}_\beta, \tilde{\mathbf{b}}_\beta\} \in \mathbb{R}^{L^p}$ . Take  $N_\gamma$  as an example, for an instance label  $l^p \in \mathbb{L}^p$  and its occupying pixels  $\mathbf{x}(l^p) = \{x | \mathbf{p}[x] = l^p\}$ , we have

$$\begin{aligned} \tilde{\mathbf{a}}_\gamma[l^p] &= (\sum_{x \in \mathbf{x}(l^p)} \tilde{\mathbf{a}}'_\gamma[x]) / |\mathbf{x}(l^p)|, \\ \tilde{\mathbf{b}}_\gamma[l^p] &= (\sum_{x \in \mathbf{x}(l^p)} \tilde{\mathbf{b}}'_\gamma[x]) / |\mathbf{x}(l^p)|. \end{aligned} \quad (4)$$

This remapping encoder, together with the main generator, forms a variational autoencoder (VAE) [22]. The remapped noise samples after perturbation are:

$$\tilde{N}_\gamma[l^p] = \tilde{\mathbf{a}}_\gamma[l^p] N_\gamma[l^p] + \tilde{\mathbf{b}}_\gamma[l^p], \quad (5)$$

where KL-Divergence loss [22] is used to enforce a same normal distribution  $\tilde{N}_\gamma \sim \mathcal{N}(0, 1)$ . These remapped noise samples are used instead of  $N_\gamma$  during modulation sampling, as described in § 3.3.

During training, the reference  $\mathbf{r}$  is exactly the ground-truth paired image. At test time, while initially sampled noises can be used to achieve random style synthesis by default, as described in § 3.3, it is also allowed to provide  $\mathbf{r}$  as instance references to control the style of the result at the instance level.

## 4. Experiments

### 4.1. Implementation Details

Our INADE generator (Figure 8) follows a similar architecture of the SPADE generator [30], but with all the SPADE layers replaced by the INADE layers. Following SPADE [30], the overall loss function consists of four loss terms: conditional adversarial loss, feature matching loss [42], perceptual loss [16] and KL-Divergence loss [22]. Details are provided in the supplementary material.

During training, by default, Adam optimizer [21] ( $\beta_1 = 0, \beta_2 = 0.9$ ) is used with fixed epoch number of 200. The learning rates for the generator and the discriminator are set to 0.0001 and 0.0004 respectively, which are gradually decreased to zero after 100 epochs. Noise  $C^0$  has 64 initialized channels, while the input noise has 256 channels, same as [30, 27]. All experiments are implemented in PyTorch [31] and conducted on TITAN XP GPUs.

### 4.2. Datasets and Metrics

**Datasets.** Experiments are conducted on five popular datasets: *Cityscapes* [3], *ADE20K* [48], *CelebAMask-HQ* [23, 17, 29], *DeepFashion* [28], and *DeepFashion2*. *DeepFashion2* is built from *DeepFashion* that each image contains two persons, which is only used for testing. More details can be found in the supplementary material.

**Quality Metrics.** To evaluate the result quality, we adopt two types of metrics following [2, 42, 30, 37]. One is the Fréchet Inception Distance (FID) [11], which measures the distance of distributions between results and real images.

The other one is semantic-segmentation-based, which evaluates the semantic segmentation accuracy on the results by comparing the predicted masks with the groundtruth layouts on both mean Intersection-over-Union (mIoU) and pixel accuracy (accu). State-of-the-art pretrained segmentation models are used for different datasets: DRN [45, 7] for *Cityscapes*, UperNet101 [43, 4] for *ADE20k*, and UNet [34, 35] for *CelebAMask-HQ*. For *DeepFashion*, we use the same UNet-based network but train the model by ourselves. For fair evaluation, we run the model for 10 times and report the average scores when noise input is required.

**Diversity Metrics.** To evaluate the diversity of the results, we adopt the LPIPS metric as proposed by [47, 33]. Similar to [52], we also adopt two metrics to measure the semantic-level diversity (mCSD and mOCD) and expand them to instance level (mISD and mOID). More details can be found

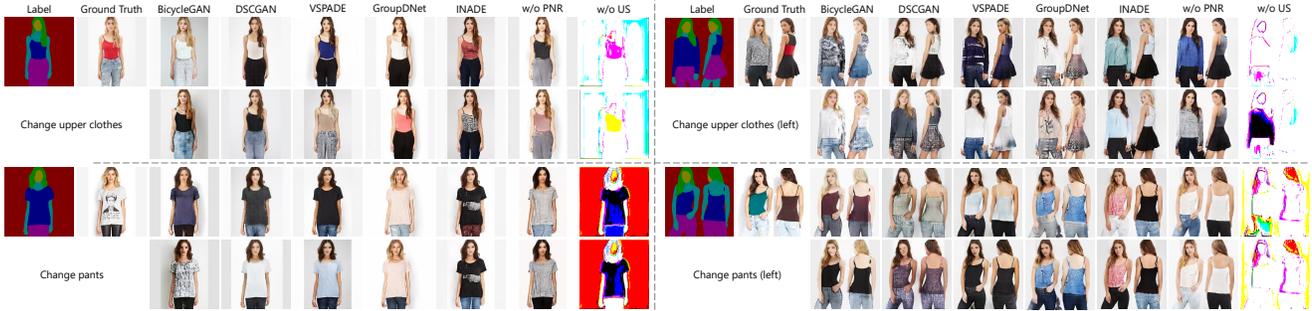


Figure 4. Visual comparison with other multimodal models and two baselines. The results on the left show the performance of class level diversity while the results on the right are for instance level diversity. The first two rows represent the results of different models when changing upper clothes while the last two rows represent the results of changing pants.

in the supplementary material.

**Subjective Metrics.** We conduct human evaluations to assess both the quality and the diversity of the methods. For quality, we ask the volunteers to select the most realistic one among the results generated by different methods on the same input. mHE (mean Human Evaluation) denotes the percentage of results being selected for each method.

For diversity, we expand the metric proposed by [52] to the instance level. A pair of results, with one random semantic class or instance manipulated, are given to volunteers. The percentage of pairs that are judged to be different in only one area represents the human evaluation, namely SHE (Semantic Human Evaluation) and IHE (Instance Human Evaluation). We invite 20 volunteers for evaluation and the evaluated number of images is 20.

### 4.3. Quantitative and Qualitative Comparisons

We compare INADE with several SOTA works, including quality-oriented (pix2pixHD [42], SPADE [30], CLADE [38, 37] and SEAN [51]) and diversity-oriented (BicycleGAN [50], DSCGAN [44], and GroupDNet [52]) methods. For a fair comparison, we directly use the pre-trained models provided by the authors when available, otherwise train the models by ourselves using the codes and settings provided by the authors. For SPADE, which has the strategy for multimodal synthesis, we train the model with that extra encoder but ignore it when testing its diversity performance (namely VSPADE). Reference images are not allowed for any of the methods during testing except for SEAN which requires the reference input.

#### 4.3.1 Multimodal Image Synthesis

Several methods that support multimodal image synthesis are compared to demonstrate the superior diversity performance of INADE. In addition, we also compare with two ablation baselines: w/o US (without Unified Sampling, §3.3) and w/o PNR (without Prior Noise Remapping, §3.4). w/o US means that the noise for each INADE layers are

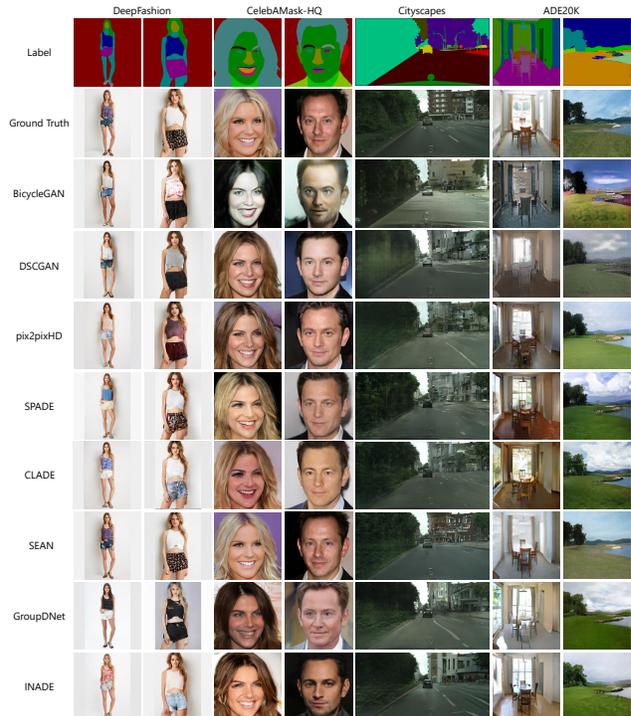


Figure 5. Qualitative comparison with the state-of-the-art semantic image synthesis methods on four datasets: DeepFashion, CelebAMask-HQ, Cityscapes and ADE20K.

sampled independently, while w/o PNR means that the PNR is not used during training. The quantitative results on the *DeepFashion* dataset are summarized in Table 1.

In general, INADE achieves superior performance regarding both quality and diversity compared to previous methods. For single-subject images in the *DeepFashion* dataset, our method exhibits better performance than BicycleGAN, DSCGAN, and VSPADE in terms of FID, and is comparable to GroupDNet. While for multiple-subject images (*DeepFashion2*), INADE shows the lowest FID.

In terms of diversity, our method achieves the best

Table 1. Comparison with other multimodal methods on diversity. mHE, SHE and IHE are aforementioned metrics.  $\uparrow$  and  $\downarrow$  represent the higher the better and the lower the better. **Bold** and underlined numbers are the best and the second best of each metric, respectively.

Methods	DeepFashion						DeepFashion2					
	FID $\downarrow$	LPIPS $\uparrow$	mCSD $\uparrow$	mOCD $\downarrow$	mHE $\uparrow$	SHE $\uparrow$	FID $\downarrow$	LPIPS $\uparrow$	mISD $\uparrow$	mOID $\downarrow$	mHE $\uparrow$	IHE $\uparrow$
BicycleGAN	31.10	<b>0.225</b>	<u>0.0465</u>	0.2014	0.0	4.5	33.46	<u>0.286</u>	<u>0.0500</u>	0.2456	0.0	2.5
DSCGAN	29.79	0.146	0.0404	0.1218	0.0	9.3	48.64	0.199	0.0433	0.1633	0.0	4.8
VSPADE	11.11	0.197	0.0450	0.1665	7.5	6.5	<u>22.29</u>	0.222	0.0390	0.1780	43.7	3.3
GroupDNet	<b>9.72</b>	0.222	0.0453	<b>0.0077</b>	<u>40.0</u>	<u>86.0</u>	22.81	0.281	0.0434	<u>0.0303</u>	8.8	<u>9.3</u>
INADE	<u>9.97</u>	<b>0.225</b>	<b>0.0511</b>	<u>0.0161</u>	<b>52.5</b>	<b>88.3</b>	<b>18.18</b>	<b>0.319</b>	<b>0.0580</b>	<b>0.0187</b>	<b>47.5</b>	<b>82.8</b>
w/o PNR	12.09	0.184	0.0289	0.0138	-	-	20.76	0.243	0.0291	0.0189	-	-
w/o US	248.33	0.624	0.0730	0.0370	-	-	265.63	0.633	0.0748	0.0296	-	-

Table 2. Comparison with SOTA methods on result quality. All the numbers are collected by running the evaluation on our machine. Here **M**, **A**, **F**, and **L** represent mIoU, accu, FID, and LPIPS, respectively. Note that the **L** score of SEAN is almost zero even with noise input.

Methods	Cityscapes				ADE20K				CelebAMask-HQ				DeepFashion			
	M $\uparrow$	A $\uparrow$	F $\downarrow$	L $\uparrow$	M $\uparrow$	A $\uparrow$	F $\downarrow$	L $\uparrow$	M $\uparrow$	A $\uparrow$	F $\downarrow$	L $\uparrow$	M $\uparrow$	A $\uparrow$	F $\downarrow$	L $\uparrow$
SPADE	<b>61.38</b>	<u>93.26</u>	51.98	0	<b>36.28</b>	<u>78.13</u>	29.79	0	75.22	<u>94.76</u>	31.40	0	<b>76.76</b>	<b>97.65</b>	11.22	0
pix2pixHD	60.50	93.06	66.04	0	27.27	72.61	45.87	0	<u>76.11</u>	<b>95.67</b>	36.95	0	73.99	97.02	15.27	0
CLADE	60.44	<b>93.42</b>	50.62	0	<u>35.43</u>	<u>77.37</u>	30.48	0	75.37	95.05	33.54	0	75.63	97.33	12.76	0
SEAN	56.22	92.28	50.43	0	32.65	76.58	<b>28.11</b>	0	75.94	95.03	<u>24.30</u>	0	<u>76.28</u>	97.46	<b>7.37</b>	0
BicycleGAN	30.47	78.26	59.87	0.122	5.33	42.68	77.49	<b>0.443</b>	65.98	89.77	35.73	<u>0.362</u>	73.09	96.75	31.10	<b>0.225</b>
DSCGAN	43.70	87.80	50.84	0.216	8.07	58.10	82.30	0.324	75.98	95.08	52.83	0.198	75.92	96.97	29.79	0.146
GroupDNet	59.20	92.78	<u>41.12</u>	0.073	26.09	73.07	39.11	0.177	<b>76.13</b>	<u>95.21</u>	29.39	0.309	76.19	<u>97.48</u>	<u>9.72</u>	0.222
INADE	<u>61.02</u>	93.16	<b>38.04</b>	<b>0.248</b>	34.96	<b>78.51</b>	<u>29.60</u>	<u>0.400</u>	74.08	94.31	<b>22.55</b>	<b>0.365</b>	76.27	97.44	9.96	<b>0.225</b>

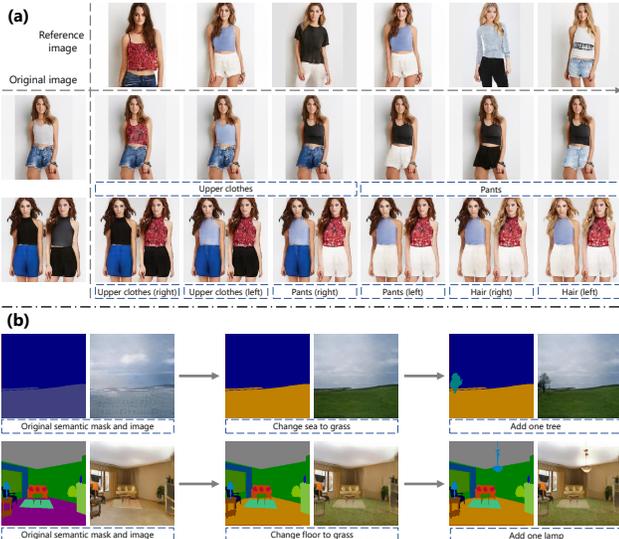


Figure 6. Exemplar applications of the proposed INADE. (a) Results of our method for reference appearance editing. From left to right, we change the appearance of part of the target image based on the reference image (the changed instance is indicated by the text in the blue dotted box). (b) Application of our method for semantic manipulation. The text in the blue dotted box indicates what is edited each time.

scores on metrics including overall measurement (LPIPS) and semantic-level/instance-level metrics (mCSD/mISD). As for mOCD, methods only support image-level diversity, such as BicycleGAN, DSCGAN, and VSPADE, produce much more unwanted changes outside the instance area.

Although both being relatively low, INADE is higher than GroupDNet, which is because GroupDNet uses group convolution to prevent premature fusion between features of different classes, while INADE uses conventional convolution for more consistent combinations. For mOID, as the only method that supports instance-wise control, INADE easily gets the best score. More analysis of mOID and mOCD can be found in the supplementary material.

As for subjective evaluations, our method also outperforms others on both semantic- and instance-level cases.

Compared to the two ablation baselines, we find that both prior noise remapping and unified sampling play indispensable roles in our method. Removing prior noise remapping (w/o PNR) leads to ambiguities during the supervised training, which seriously affects the quality and diversity of synthesized results. Independently sampling (w/o US) for each normalization destroys the consistency of information and significantly degenerates the generation result.

The qualitative comparisons are shown in Figure 4. Although all methods support multimodal synthesis, the quality by BicycleGAN and DSCGAN is not satisfactory. VSPADE achieves good visual quality, but does not support semantic- or instance-level control. GroupDNet is capable of changing the appearance of a specific semantic class (results on the left), but tends to generate identical style for different cloth instances (results on the right). On the contrary, INADE supports both fine-grained multimodal controls with high visual fidelity. As for the ablation baselines, we notice that removing PNR significantly decreases the quality, while the whole task fails without US.

Table 3. Comparison with other semantic image synthesis methods on model complexity and efficiency. All the numbers are collected by running the evaluation on Titan XP. Here **P** and **T** denote the number of generator parameters and inference run time, respectively.

Methods	Cityscapes			ADE20K			CelebAMask-HQ			DeepFashion		
	<b>P</b> (M)	FLOPs (G)	<b>T</b> (s)	<b>P</b> (M)	FLOPs (G)	<b>T</b> (s)	<b>P</b> (M)	FLOPs (G)	<b>T</b> (s)	<b>P</b> (M)	FLOPs (G)	<b>T</b> (s)
SPADE	93.05	281.54	0.065	96.50	181.30	0.042	92.54	141.32	0.035	92.21	137.99	0.032
pix2pixHD	182.53	151.32	0.038	182.90	99.30	0.041	182.47	72.17	0.023	182.44	69.91	0.020
CLADE	67.90	75.54	0.035	71.40	42.20	0.024	67.32	42.15	0.022	66.98	42.15	0.019
SEAN	330.41	681.75	0.507	-	-	-	266.90	346.27	0.165	223.23	342.89	0.135
BicycleGAN	54.80	18.40	0.011	54.80	18.40	0.006	54.80	18.40	0.006	54.80	18.40	0.006
DSCGAN	54.00	18.14	0.018	54.00	18.14	0.010	54.00	18.14	0.010	54.00	18.14	0.010
GroupDNet	76.50	463.61	0.224	68.33	383.00	0.088	145.29	225.53	0.090	96.32	291.61	0.062
INADE	77.39	75.25	0.048	90.89	42.19	0.035	85.12	42.18	0.030	84.63	42.92	0.026

### 4.3.2 Semantic Image Synthesis

The quantitative comparisons against semantic image synthesis methods are summarized in Table 2.

Compared to methods that don’t support multimodal (i.e. 0 LPIPS), especially SEAN which has additional reference image input, INADE has an advantage or near the best on almost all metrics. It seems that SPADE has slight advantage in segmentation metrics, but INADE still shows its overall superiority when considering FID score and visual results.

Compared to existing multimodal methods, INADE leads all metrics. In terms of quality (e.g. mIoU, acc, and FID), BicycleGAN and DSCGAN are much lower than ours on all datasets. GroupDNet achieves similar or slightly better performance on *CelebAMask-HQ* and *DeepFashion* datasets, but has a significant gap on more complicated scenes such as *Cityscapes* and *ADE20K*. This demonstrates the superiority of the proposed method on synthesizing complex scenes. The LPIPS score shows that all these methods are able to generate multimodal images to some extent. BicycleGAN gets higher scores than ours on some datasets, but is not able to do high-quality synthesis. GroupDNet shows good performance on person-related tasks, but falls into strong bias when dealing with complex scenes which greatly restricts its performance. Therefore, considering both quality and diversity, our method achieves the best overall performance.

Qualitative comparisons on these four datasets are shown in Figure 5. In general, the images generated by INADE are more realistic than others on various datasets, which is consistent with the quantitative results.

### 4.3.3 Computational and Model Complexity

In this section, we analyze the computational and model overhead of different methods. The quantitative results (generator networks only) are summarized in Table 3.

BicycleGAN and DSCGAN share a similar small network architecture with the least parameters, FLOPs (floating-point operations per second), and run-time cost. However, the quality of the synthesized images is far from satisfactory. CLADE seems to get a good trade-off between

performance and efficiency, but still falls short of INADE in terms of overall performance and functionality. Compared to all other methods, INADE achieves the smallest network (parameters and FLOPs), as well as one of the fastest run-time performance. Specifically, compared with GroupDNet, the only method that can achieve semantic-level diversity, our method provides control over both semantic and instance levels with much less overhead, introducing 82% ~ 89% fewer FLOPs and 59% ~ 79% less inference time compared with GroupDNet.

## 4.4. Applications

Thanks to its superior capability of controllable diverse synthesis, INADE can be used in many image editing applications. Here we show two examples as follows.

**Reference appearance editing.** With the noise remapping mechanism described in § 3.4, we can extract the instance-wise style from an arbitrary reference image. This makes it possible for INADE to perform reference-based editing to different parts of an image at the instance level. As shown in Figure 6 (a), we can change the appearance of hairstyles, tops, and pants to match the reference.

**Semantic manipulation.** Similar to most existing semantic image synthesis methods, INADE also supports semantic manipulation. We show some examples in Figure 6 (b), such as changing the semantic class to an object, or insert a new semantic object into the image. And more creative editing results can be achieved by modifying the semantic mask and the instance map.

## 5. Conclusion

In this paper, we focus on multimodal image synthesis and propose a novel diverse semantic image synthesis method based on instance-aware conditional normalization. Different from previous works, we learn the class-wise probability distributions and perform instance-wise stochastic sampling to generate the per-instance modulation parameters. Our method improves the network’s ability to model semantic categories and make it easier to synthesize diverse images at semantic- or instance-level without scarifying the visual fidelity.

## References

- [1] Amjad Almahairi, Sai Rajeswar, Alessandro Sordani, Philip Bachman, and Aaron Courville. Augmented cycleGAN: Learning many-to-many mappings from unpaired data. *arXiv preprint arXiv:1802.10151*, 2018. **2**
- [2] Qifeng Chen and Vladlen Koltun. Photographic image synthesis with cascaded refinement networks. In *Proceedings of the IEEE international conference on computer vision*, pages 1511–1520, 2017. **5**
- [3] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. **2, 5, 11**
- [4] CSAILVision. Pytorch implementation for semantic segmentation/scene parsing on mit ade20k dataset. <https://github.com/CSAILVision/semantic-segmentation-pytorch.git>, 2019. **5**
- [5] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. *arXiv preprint arXiv:1610.07629*, 2016. **3**
- [6] Aysegül Dundar, Karan Sapra, Guilin Liu, Andrew Tao, and Bryan Catanzaro. Panoptic-based image synthesis. *arXiv preprint arXiv:2004.10289*, 2020. **2, 3**
- [7] fyu. Dilated residual networks. <https://github.com/fyu/drn.git>, 2019. **5**
- [8] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. **1**
- [9] Shuyang Gu, Jianmin Bao, Hao Yang, Dong Chen, Fang Wen, and Lu Yuan. Mask-guided portrait editing with conditional GANs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3436–3445, 2019. **2**
- [10] Adam W Harley, Konstantinos G Derpanis, and Iasonas Kokkinos. Segmentation-aware convolutional networks using local attention masks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5038–5047, 2017. **5**
- [11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637, 2017. **5**
- [12] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1501–1510, 2017. **3**
- [13] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 172–189, 2018. **2**
- [14] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456, 2015. **3**
- [15] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. **1, 2, 11, 13**
- [16] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. **5**
- [17] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. **2, 5, 11**
- [18] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019*, pages 4401–4410, 2019. **1**
- [19] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of styleGAN. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020. **1**
- [20] Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, and Jiwon Kim. Learning to discover cross-domain relations with generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1857–1865. JMLR. org, 2017. **2**
- [21] D Kingma and J Ba. Adam: A method for stochastic optimization in: Proceedings of international conference on learning representations. 2015. **5**
- [22] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. **5**
- [23] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. MaskGAN: Towards diverse and interactive facial image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5549–5558, 2020. **2, 5, 11**
- [24] Hsin-Ying Lee, Hung-Yu Tseng, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. Diverse image-to-image translation via disentangled representations. In *Proceedings of the European conference on computer vision (ECCV)*, pages 35–51, 2018. **2**
- [25] Guilin Liu, Fitsum A Reda, Kevin J Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. Image inpainting for irregular holes using partial convolutions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 85–100, 2018. **5**
- [26] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in neural information processing systems*, pages 700–708, 2017. **2**
- [27] Xihui Liu, Guojun Yin, Jing Shao, Xiaogang Wang, et al. Learning to predict layout-to-image conditional convolutions for semantic image synthesis. In *Advances in Neural Information Processing Systems*, pages 568–578, 2019. **2, 5**
- [28] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. DeepFashion: Powering robust clothes recognition

- and retrieval with rich annotations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1096–1104, 2016. **2, 5, 11**
- [29] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015. **2, 5, 11**
- [30] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2337–2346, 2019. **1, 2, 3, 5, 6**
- [31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019. **5**
- [32] Xiaojuan Qi, Qifeng Chen, Jiaya Jia, and Vladlen Koltun. Semi-parametric image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8808–8816, 2018. **2, 11**
- [33] richzhang. Perceptualsimilarity. <https://github.com/richzhang/PerceptualSimilarity.git>, 2020. **5, 11**
- [34] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. **5, 12**
- [35] switchablenorms. Celebamask-hq. <https://github.com/switchablenorms/CelebAMask-HQ.git>, 2020. **5**
- [36] Zhentao Tan, Menglei Chai, Dongdong Chen, Jing Liao, Qi Chu, Lu Yuan, Sergey Tulyakov, and Nenghai Yu. Michigan: multi-input-conditioned hair image generation for portrait editing. *ACM Transactions on Graphics (TOG)*, 39(4):95–1, 2020. **1**
- [37] Zhentao Tan, Dongdong Chen, Qi Chu, Menglei Chai, Jing Liao, Mingming He, Lu Yuan, Gang Hua, and Nenghai Yu. Semantic image synthesis via efficient class-adaptive normalization. *arXiv preprint arXiv:2012.04644*, 2020. **2, 3, 4, 5, 6**
- [38] Zhentao Tan, Dongdong Chen, Qi Chu, Menglei Chai, Jing Liao, Mingming He, Lu Yuan, and Nenghai Yu. Rethinking spatially-adaptive normalization. *arXiv preprint arXiv:2004.02867*, 2020. **2, 6**
- [39] Hao Tang, Song Bai, and Nicu Sebe. Dual attention gans for semantic image synthesis. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 1994–2002, 2020. **2**
- [40] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016. **1**
- [41] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. *arXiv preprint arXiv:1808.06601*, 2018. **11**
- [42] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018. **1, 2, 3, 5, 6, 11**
- [43] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 418–434, 2018. **5**
- [44] Dingdong Yang, Seunghoon Hong, Yunseok Jang, Tianchen Zhao, and Honglak Lee. Diversity-sensitive conditional generative adversarial networks. In *International Conference on Learning Representations*, 2018. **1, 2, 6**
- [45] Fisher Yu, Vladlen Koltun, and Thomas Funkhouser. Dilated residual networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 472–480, 2017. **5**
- [46] Pan Zhang, Bo Zhang, Dong Chen, Lu Yuan, and Fang Wen. Cross-domain correspondence learning for exemplar-based image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5143–5153, 2020. **2**
- [47] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018. **5, 11**
- [48] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017. **2, 5, 11**
- [49] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. **2**
- [50] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In *Advances in neural information processing systems*, pages 465–476, 2017. **1, 2, 6**
- [51] Peihao Zhu, Rameen Abdal, Yipeng Qin, and Peter Wonka. Sean: Image synthesis with semantic region-adaptive normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5104–5113, 2020. **1, 2, 6**
- [52] Zhen Zhu, Zhiliang Xu, Ansheng You, and Xiang Bai. Semantically multi-modal image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5467–5476, 2020. **1, 2, 3, 5, 6, 11, 12**

## 6. Additional Implementation details

**Network architectures.** Here we give detailed network designs for each part. Figure 7 shows the architecture of our encoder. We use instance partial convolution and instance average pooling to get the parameters of each instance independently. The architecture of generator network is shown in Figure 8. The synthesis process starts with a random noise and goes through a series of the proposed INADE ResBIKs. Since the training is carried out on multiple GPUs, the batch normalization layer in INADE adopts the synchronous version. We use a multi-scale PathGAN [15] based discriminator whose architecture is shown in Figure 9.

**Loss function.** The loss function we adopted consists of four components:

*Conditional adversarial loss.* Let  $\mathcal{E}$  be the prior noise remapping,  $G$  be the INADE generator,  $D$  be the discriminator,  $\mathbf{m}$  be a given semantic mask,  $\mathbf{o}$  and  $\mathbf{p}$  be the corresponding image and instance map. The conditional adversarial loss built with hinge loss is formulated as:

$$\mathcal{L}_{GAN}(\mathcal{E}, G, D) = \mathbb{E}[\max(0, 1 - D(\mathbf{o}, \mathbf{m}, \mathbf{p}))] + \mathbb{E}[\max(0, 1 + D(G(\mathcal{E}(\mathbf{o}, \mathbf{p}), \mathbf{m}, \mathbf{p}), \mathbf{m}, \mathbf{p}))]. \quad (6)$$

*Feature matching loss.* Let  $D_i$  and  $N_i$  be the output feature maps and the number of elements of the  $i$ -the layer of  $D$  respectively,  $S_D$  and  $E_D$  be the start number of layer for loss calculation and total number layers in  $D$  respectively. The feature matching loss is denoted as:

$$\mathcal{L}_F = \mathbb{E} \sum_{i=S_D}^{E_D} \frac{1}{N_i} [\|D_i(\mathbf{o}, \mathbf{m}, \mathbf{p}) - D_i(G(\mathcal{E}(\mathbf{o}, \mathbf{p}), \mathbf{m}, \mathbf{p}), \mathbf{m}, \mathbf{p}))\|_1]. \quad (7)$$

To reduce the ambiguity, we only use high-level features and set  $S_D$  to 3.

*Perceptual loss.* Let  $V_i$  and  $M_i$  be the output feature maps and the number of elements of the  $i$ -the layer of VGG network respectively,  $S_V$  and  $E_V$  be the start number of layer for loss calculation and total number layers in VGG network respectively. The perceptual loss is denoted as:

$$\mathcal{L}_P = \mathbb{E} \sum_{i=S_V}^{E_V} \frac{1}{M_i} [\|V_i(\mathbf{o}) - V_i(G(\mathcal{E}(\mathbf{o}, \mathbf{p}), \mathbf{m}, \mathbf{p}))\|_1]. \quad (8)$$

Similar to feature matching loss, we only use high-level features and set  $S_D$  to 3.

*KL-Divergence loss.* Let  $q_\beta(z|\mathbf{o}, \mathbf{p})$  and  $q_\gamma(z|\mathbf{o}, \mathbf{p})$  be the variational distribution of  $N_\gamma$  and  $N_\beta$  respectively.  $p(z)$  be a standard Gaussian distribution. The KL-Divergence loss is denoted as:

$$\mathcal{L}_{KL} = 0.5 \times (\mathcal{D}(q_\beta(z|\mathbf{o}, \mathbf{p})\|p(z)) + \mathcal{D}(q_\gamma(z|\mathbf{o}, \mathbf{p})\|p(z))). \quad (9)$$

The overall loss is made up of the above-mentioned loss terms as:

$$\min_{\mathcal{E}, G} (\max_D (\mathcal{L}_{GAN}) + \lambda_1 \mathcal{L}_F + \lambda_2 \mathcal{L}_P + \lambda_3 \mathcal{L}_{KL}), \quad (10)$$

Following SPADE, We set  $\lambda_1 = 10$ ,  $\lambda_2 = 10$ ,  $\lambda_3 = 0.05$ .

## 7. Details of Datasets

The details about each dataset are described as follows:

- *Cityscapes* dataset [3] is a widely used dataset for semantic image synthesis [41, 32, 42]. The high-resolution images with fine semantic and instance annotations are taken from street scenes of German cities. There are 2,975 training images and 500 validation images. The number of annotated semantic classes is 35.
- *ADE20K* dataset [48] consists of 25,210 images (20,210 for training, 2,000 for validation and 3,000 for testing). The images in *ADE20K* dataset cover a wide range of scenes and object categories, including a total of 150 object and stuff classes.
- *CelebAMask-HQ* dataset [23, 17, 29] is based on CelebAHQ face image dataset. It contains of 28,000 training images and 2,000 validation images with 19 different semantic classes.
- *DeepFashion* dataset [28] contains 52,712 person images with fashion clothes. We use the processed dataset provided by GroupDNet [52] which consists of 30,000 training images and 2,247 validation images. There are 8 different semantic classes.
- *DeepFashion2* dataset is built from DeepFashion. We combine two adjacent images to generate the images containing two persons. The new semantic mask and the instance map are also derived from the corresponding two semantic masks. This dataset is only used to evaluate the performance of models trained on DeepFashion dataset in terms of instance level diversity.

In these datasets, *Cityscapes* and *DeepFashion2* have semantic and instance annotations, while the rest have only semantic annotations. In our experiment, the resolution of images is  $256 \times 256$  except that *Cityscapes* dataset is  $256 \times 512$ .

## 8. Details of Diversity Metrics

We adopt the LPIPS [47, 33] to evaluate the overall diversity of the results. Specifically, we generate 10 groups of images or evaluation with randomly sampled noise, and calculate the diversity score between 2 random groups at a time. A total of 10 scores are calculated, and we measure the mean of these scores to reduce the potential fluctuation caused by random sampling.

To evaluate the instance-level diversity, we expand the metrics proposed by [52], called mean *Instance-Specific Diversity* (mISD) and mean *Other-Instances Diversity* (mOID), which represent the degree of change inside and outside the instance region when being manipulated. Specially, we generate several images by changing sampled noise for specified instance while keeping the noise for others unchanged. Then, the similarity inside and outside the instance region between these images makes up the mISD and mOID metrics. For datasets which have no instance annotations, these metrics degenerate to semantic level (mean *Class-Specific Diversity* (mCSD) and mean *Other-Classes Diversity* (mOCD)) which are the same with [52]. A high diversity inside the instance area (high mISD), as well as a low outside diversity (low mOID), are desired.

### 9. Additional ablation study

Here we give the additional ablation study for  $C^0$  which represents the length of the initial sampling. Intuitively, the longer the sampling length is, the higher the diversity of the synthesized image will be. We conduct experiments on the *Cityscapes* and *CelebAMask-HQ* datasets, which include complex street scenes and delicate facial images. As summarized in Table 4, we compare the default setting ( $C^o = 64$ ) with two variant settings: a shorter sampling length ( $C^o = 8$ , INADE-8) and a longer sampling length ( $C^o = 128$ , INADE-128). We find that INADE-8 shows the lower LPIPS score than INADE, while INADE-128 correspondingly gets the highest score in this metric. And the model with the default setting (INADE) gets the best scores in terms of quality metrics. In our understanding, a short sampling length (e.g. 8) may limit the information capacity, thus reducing the generation quality (low scores of mIoU, acc and FID) and diversity (low score of LPIPS). In contrast, a longer sampling length (e.g. 128) can increase the diversity of the synthesized image (high score of LPIPS), but also increases the difficulty of high-quality image generation (low scores of mIoU, acc and FID).

In terms of model parameters, FLOPs and run time, INADE-8 is best, but the advantage is not obvious compared with INADE and INADE-128. Based on the above results, we set  $C^o = 64$  on different datasets by default.

### 10. Additional results

In Figure 10, we show more multi-modal qualitative results on different datasets that only change one specified class or instance. The conclusions are basically the same as we mention in the main submission. BicycleGAN, DSCGAN and VSPADE show the global style controllability, GroupDNet expands it to semantic level, while the synthesis results of our method can be controlled at both the semantic level and instance level. We notice that in some results,

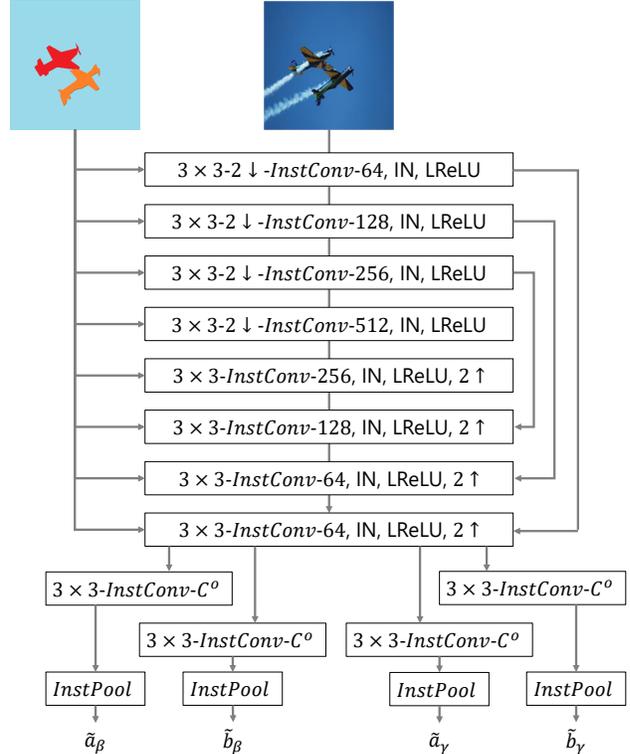


Figure 7. Architecture of our encoder network. We use UNet [34] based network to extract the features with the same resolution of input image, and then obtain the  $\{\tilde{a}_\gamma, \tilde{b}_\gamma, \tilde{a}_\beta, \tilde{b}_\beta\}$  through independent instance partial convolution (InstConv) and instance average pooling (InstPool).

when we change one part, other parts slightly change as well, which is also mentioned in GroupDNet [52]. In fact, this is reasonable in some cases to increase the generation fidelity. For example, as shown in Figure 10 (h), the lighting often changes with the sky, if the appearance of the grass is totally unchanged, the final generated image will look unnatural to some extent. Therefore, though the metric mOCD (or mOID) may be a good indication of semantic/instance-level controllability, a slightly high mOCD or mOID do not represent worse quality. In other words, we do not expect them to be zero in real applications.

In Figure 11, Figure 12, Figure 13, Figure 14, we further show additional qualitative comparison results between the proposed INADE and other methods on the *DeepFashion*, *Cityscapes*, *ADE20K* and *CelebAMask-HQ* datasets. These results show that the images quality of INADE is better than or at least comparable to existing methods.

Table 4. Comparison of INADE with different  $C^o$  on the Cityscapes and CelebAMask-HQ daasets. **P**, **F** and **T** represent the generator parameters, FLOPs and run time respectively.

Methods	Cityscapes							CelebAMask-HQ						
	mIoU	acc	FID	LPIPS	<b>P</b> (M)	<b>F</b> (G)	<b>T</b> (s)	mIoU	acc	FID	LPIPS	<b>P</b> (M)	<b>F</b> (G)	<b>T</b> (s)
INADE-64 (default)	<b>61.02</b>	<b>93.16</b>	<b>38.04</b>	0.248	77.39	75.26	0.0486	<b>74.08</b>	<b>94.31</b>	<b>22.55</b>	0.365	85.12	42.18	0.0298
INADE-8	60.25	93.07	38.68	0.220	<b>76.78</b>	<b>75.23</b>	<b>0.0482</b>	73.26	<b>94.31</b>	24.58	0.350	<b>84.50</b>	<b>42.16</b>	<b>0.0295</b>
INADE-128	59.57	92.68	39.30	<b>0.315</b>	78.10	75.28	0.0497	73.48	94.28	24.88	<b>0.366</b>	85.82	42.20	0.0306

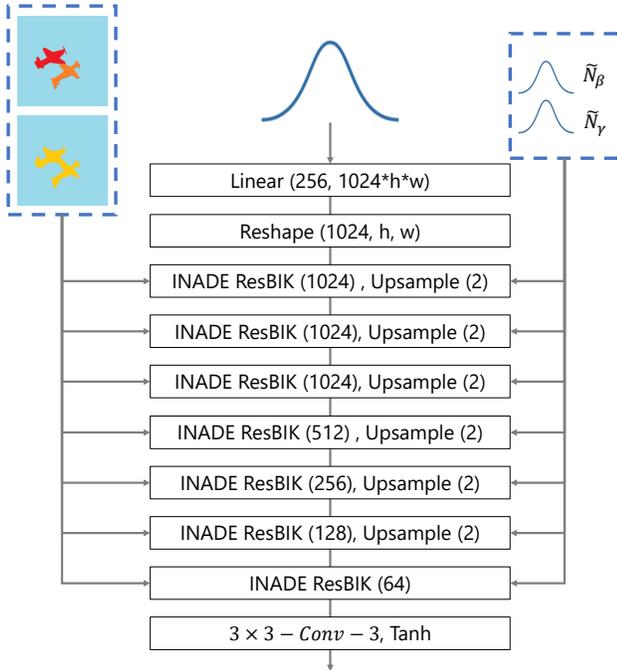


Figure 8. Architecture of our generator network. It consists of a linear transform layer, six INADE ResBIKs with upsampling and a final classification convolution layer. The upsampling operation on the second INADE ResBIK is removed if the resolution of generated images is  $256 \times 512$ . The initial noise ( $\tilde{N}_\gamma, \tilde{N}_\beta$ ) will be translated through a linear transformation mapping before fed to INADE ResBIKs.

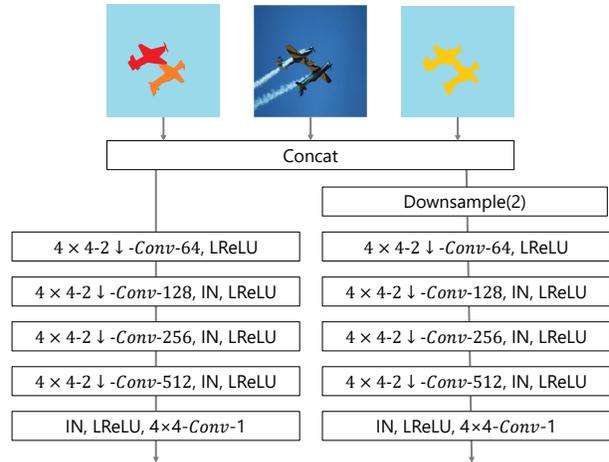


Figure 9. The discriminator of our method is based on the PatchGAN [15]. It takes the concatenation the segmentation map, instance map and the image as input.

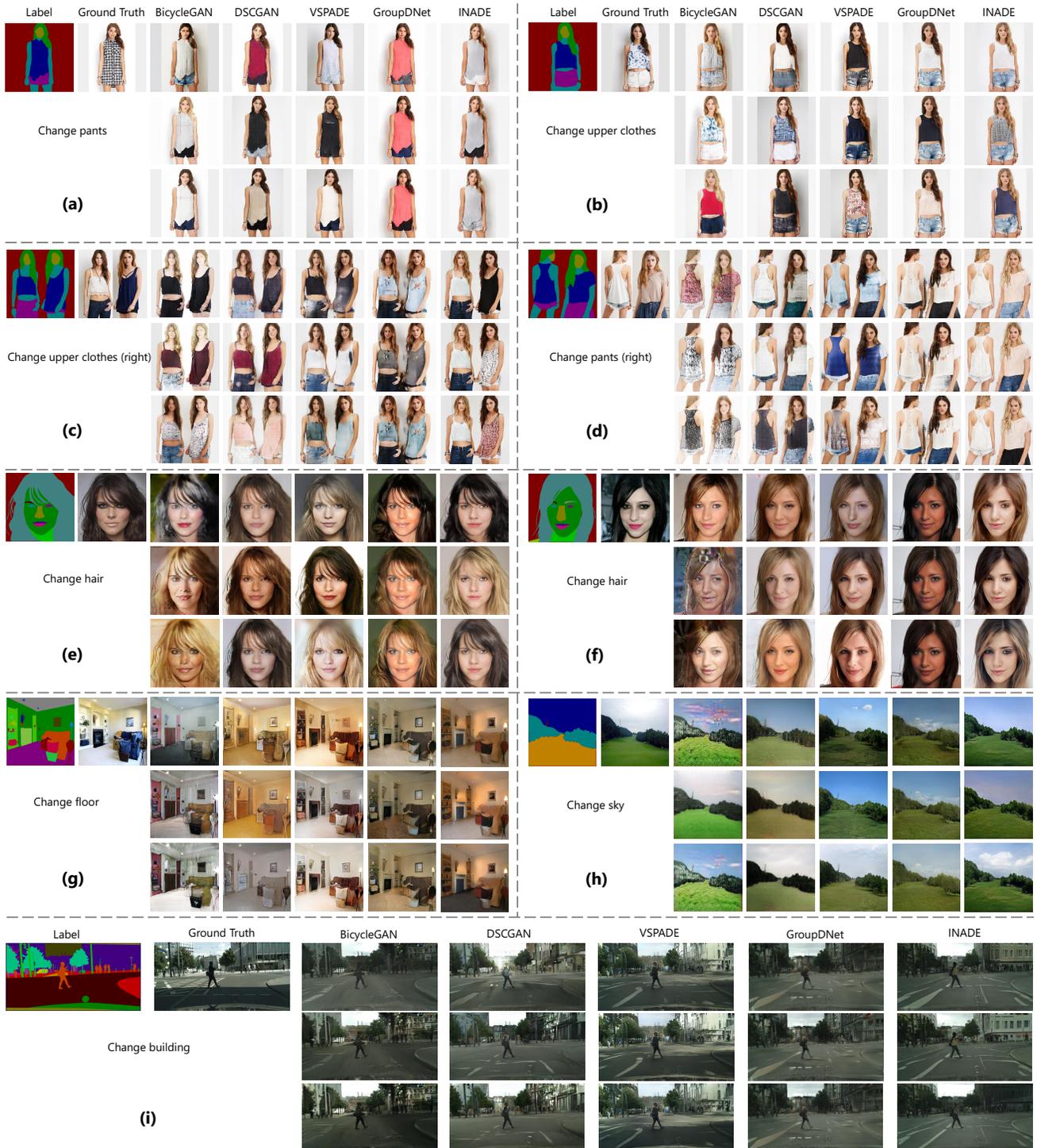


Figure 10. Multi-modal comparison of our INADE with previous state-of-the-art methods on DeepFashion (a-b), DeepFashion2 (c-d), CelebAMask-HQ (e-f), ADE20K (g-h) and Cityscapes (i) datasets.



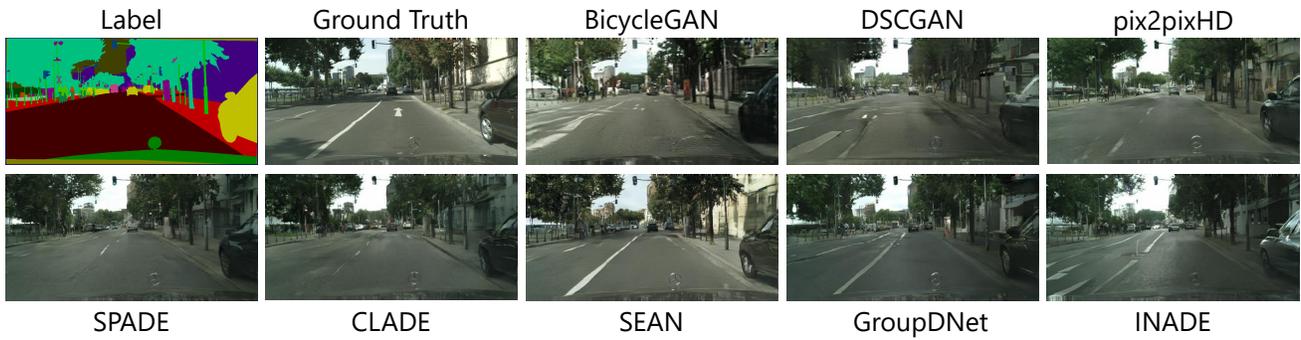
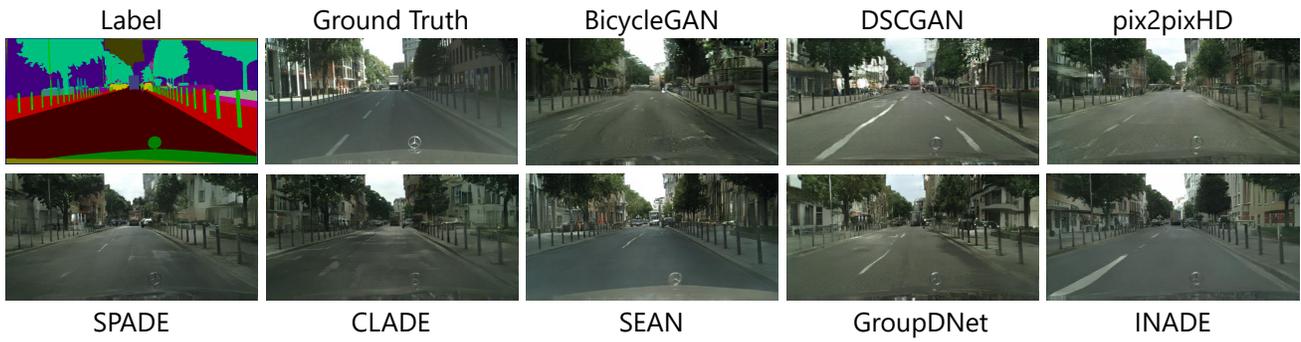
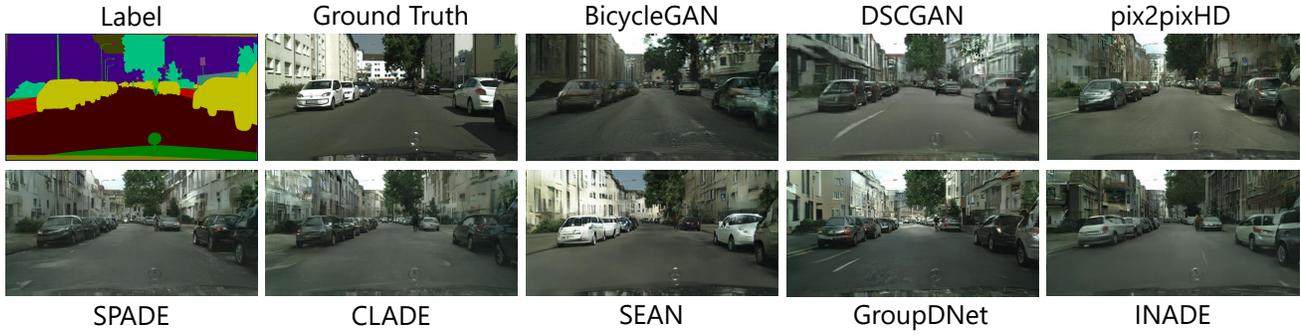


Figure 12. Qualitative comparison of our INADE with previous state-of-the-art methods on the Cityscapes dataset.



Figure 13. Qualitative comparison of our INADE with previous state-of-the-art methods on the ADE20K dataset.

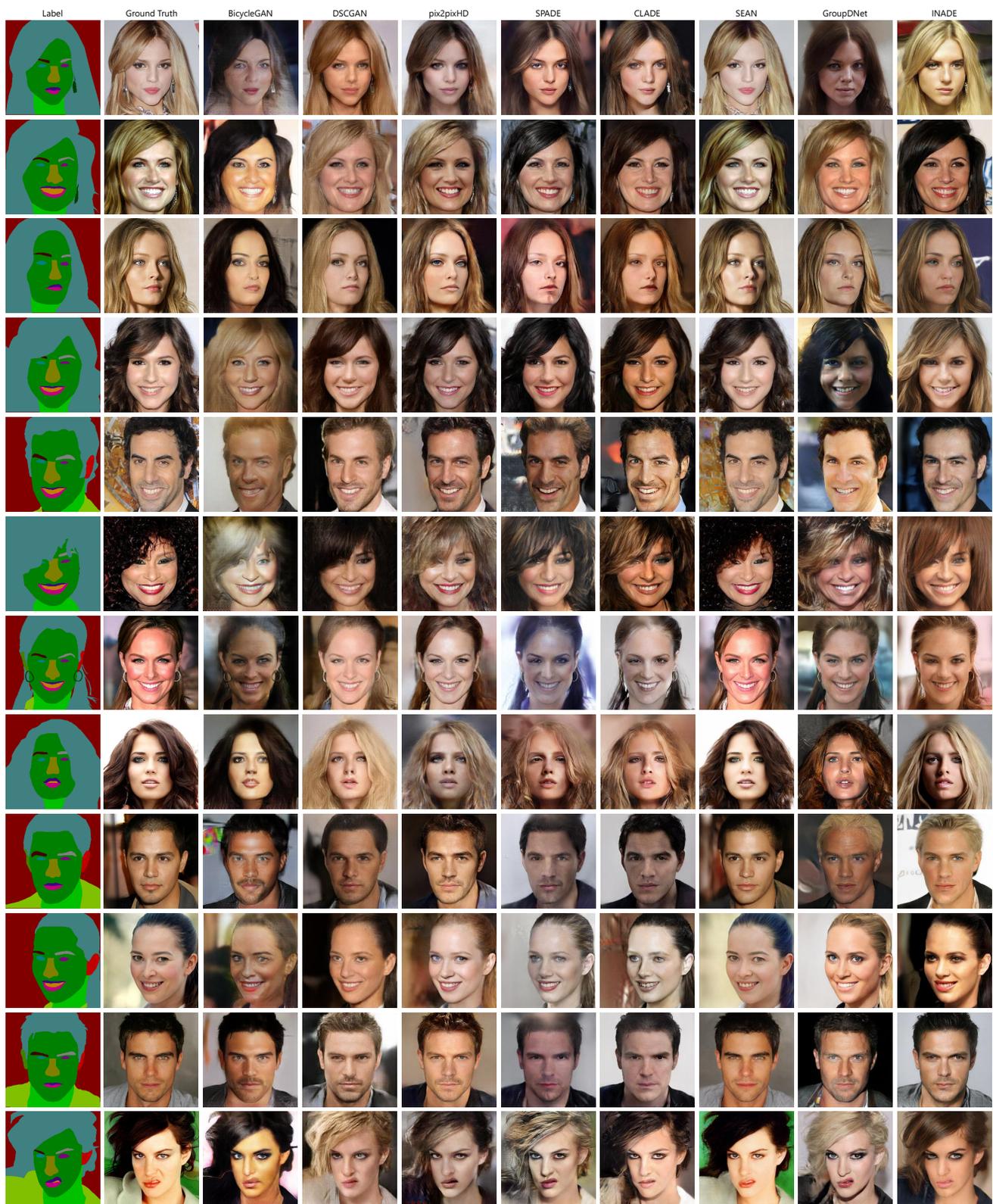


Figure 14. Qualitative comparison of our INADE with previous state-of-the-art methods on the CelebAMask-HQ dataset.