# Attribute Guided Unpaired Image-to-Image Translation
# with Semi-supervised Learning

Xinyang Li[1#], Jie Hu[1#], Shengchuan Zhang[1], Xiaopeng Hong[2],

Qixiang Ye[3], Chenglin Wu[4], and Rongrong Ji[1*]

[1]Xiamen University, [2]Xi'an Jiaotong University, [3]University of Chinese Academy of Sciences, [4]Fuzhi.ai

{Imlixinyang,hujie.cpp,ether.wcl}@gmail.com, {rrji, zsc_2016}@xmu.edu.cn

hongxiaopeng@mail.xjtu.edu.cn, qxye@ucas.ac.cn

## Abstract

*Unpaired Image-to-Image Translation (UIT) focuses on translating images among different domains by using unpaired data, which has received increasing research focus due to its practical usage. However, existing UIT schemes defect in the need of supervised training, as well as the lack of encoding domain information. In this paper, we propose an Attribute Guided UIT model termed AGUIT to tackle these two challenges. AGUIT considers multimodal and multi-domain tasks of UIT jointly with a novel semi-supervised setting, which also merits in representation disentanglement and fine control of outputs. Especially, AGUIT benefits from two-fold: (1) It adopts a novel semisupervised learning process by translating attributes of labeled data to unlabeled data, and then reconstructing the unlabeled data by a cycle consistency operation. (2) It decomposes image representation into domain-invariant content code and domain-specific style code. The redesigned style code embeds image style into two variables drawn from standard Gaussian distribution and the distribution of domain label, which facilitates the fine control of translation due to the continuity of both variables. Finally, we introduce a new challenge, i.e., disentangled transfer, for UIT models, which adopts the disentangled representation to translate data less related with the training set. Extensive experiments demonstrate the capacity of AGUIT over existing state-of-the-art models.*

## 1. Introduction

Image-to-image translation aims to learn the mapping between images among different domains, which has drawn increasing research attention. Many computer vision tasks can be modeled as an image-to-image translation problem,

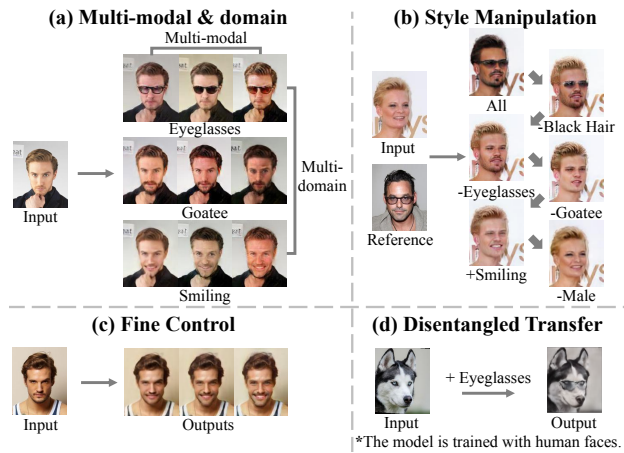#Contributed Equally.

*Corresponding Author.

Figure 1. Examples of UIT tasks accomplished by AGUIT. (a) Multi-modal and multi-domain translation are achieved jointly. (b) Style code of reference image can be manipulated while retaining the content of input image. (c) The outputs can be fine controlled. (d) An example of disentangled transfer.

such as colorization [38], super resolution [21, 36], semantic synthesis [4] and domain adaption [14]. Earliest works in image-to-image translation can be referred to translate paired images (*i.e.*, Pix2Pix [17]), in which every image in the source domain should have the corresponding paired image in the target domain for training. However, it is expensive to collect paired images in practice. To overcome this limitation, several methods were proposed for unpaired image-to-image translation (UIT), such as CycleGAN [40], DualGAN [37] and UNIT [24].

UIT involves two basic tasks, *i.e.*, multi-domain translation and multi-modal translation. The former aims to use a single model to achieve translations across multiple domains. The latter aims to generate diversified outputs while retaining the content information. For multi-domain translation, StarGAN [8] was proposed to input the labels together with images into the model to guide the translation.

| Merits | Unpaired | Multi-modal | Multi-domain | Single Model | Fine Control | Disentanglement | Semi Supervised | Disentangled Transfer |
|---|---|---|---|---|---|---|---|---|
| Pix2Pix [17] | - | - | - | - | - | - | - | - |
| BicycleGAN [41] | - | √ | - | - | - | - | - | - |
| CycleGAN [40] | √ | - | - | - | - | - | - | - |
| MUNIT [16] / DRIT [22] | √ | √ | - | - | - | - | - | - |
| StarGAN [8] | √ | - | √ | √ | - | - | - | - |
| GANimation [32] | √ | - | √ | √ | √ | - | - | - |
| SMIT [33] | √ | √ | √ | √ | - | - | - | - |
| **AGUIT** | √ | √ | √ | √ | √ | √ | √ | √ |

Table 1. The merits of AGUIT compared with the existing state-of-the-art image-to-image translation models.

However, due to the fixed and discrete labels, it cannot generate diverse outputs. GANimation [32] describes the labels in a continuous manifold, which however severely relies on the Action Units annotations. For multi-modal translation, MUNIT [16] and DRIT [22] were proposed, in which the style code is drawn from the Gaussian distribution to enrich the outputs. However, these two methods need different models for different translations. Recently, SMIT [33] is proposed to achieve these two tasks jointly, in which the images with labels and additional style noises are input to the generator. However, it cannot control the result finely since the domain information is not encoded. Overall, these state-of-the-art UIT models are defect in the need of supervised training, as well as a suitable and explicit encoding for domain information.

We argue that an ideal UIT model should have the following merits: First, the unlabeled data should be incorporated into training process to achieve *semi-supervised learning*, so as to reduce the requirement of expensive label annotations. Second, the domain information (or attributes) of images should be explicitly learned to *finely control* the result, so as to enhance the translation flexibility. Third, the learned representation should be *disentangled*, so the inputs can be translated by specific semantic information to improve the translation interpretability.

In this paper, we achieve the above goals in a unified framework by proposing an Attribute Guided Unpaired Image-to-image Translation (AGUIT) model. Two innovative designs are presented: (1) It adopts a novel semi-supervised learning process by translating attributes of labeled data to unlabeled data, and then reconstructing the unlabeled data with a cycle consistency operation. (2) It decomposes image representation into domain-invariant content code and domain-specific style code. The style code embeds image style into two variables drawn from standard Gaussian distribution and domain label distribution, which enhances the controllability of translation due to the continuity of both distributions. Particularly, the standard Gaussian distribution encourages unsupervised disentanglement, since its every dimension is independent with each other. And the domain label distribution is disentangled naturally and forced to be continuous by our training scheme. Finally, we introduce a *disentangled transfer*, which is a new challenge for UIT to adopt the disentangled representation to translate data less related with training set. The examples of UIT tasks accomplished by AGUIT are shown in Fig. 1.

Extensive experiments demonstrate the capacity of AGUIT over the state-of-the-art image-to-image translation models. First, AGUIT performs well on the basic tasks (*i.e.*, multi-domain translation and multi-modal translation) of UIT. Second, qualitative and quantitative evaluations reveal the benefits of our semi-supervised scheme. Third, AGUIT works well for high-level UIT tasks such as style manipulation and fine control. Finally, AGUIT carries out disentangled transfer, which is a new task introduced to image-to-image translation. The merits of AGUIT are summarized as Tab. 1. Our contributions are listed as below:

- The proposed AGUIT, which employs a novel semi-supervised learning process and a novel style code for translation. To our best knowledge, AGUIT is the first model to consider multi-modal translation and multi-domain translation with a semi-supervised setting, and achieves disentanglement of representation as well as fine control of outputs for UIT.

- We introduce a new translation task termed disentangled transfer, in which the disentangled representation is adopted to translate data less related with training set.

- Extensive experiments demonstrate the capacity of AGUIT over the existing state-of-the-art image-to-image translation models.

The rest of this paper is organized as follows. Sec. 2 reviews the related work. The proposed AGUIT is introduced in Sec. 3. Qualitative and quantitative experiments are given in Sec. 4. Finally, we conclude this work in Sec. 5.

## 2. Related Work

**Generative Adversarial Network.** GAN [12] has achieved remarkable results. In the training stage of GAN, the discriminator tries to distinguish the outputs of generator and the real distribution. On the contrary, the generator tries to fool the discriminator. After training, the generator can produce outputs which are similar to the real samples. In our model, GAN is used to align the labels for inputs, constrain the content code to common space, and make the style code and content code be separable.

**Semi-supervised GANs.** Several recent works leveraged GANs for semi-supervised learning of classification
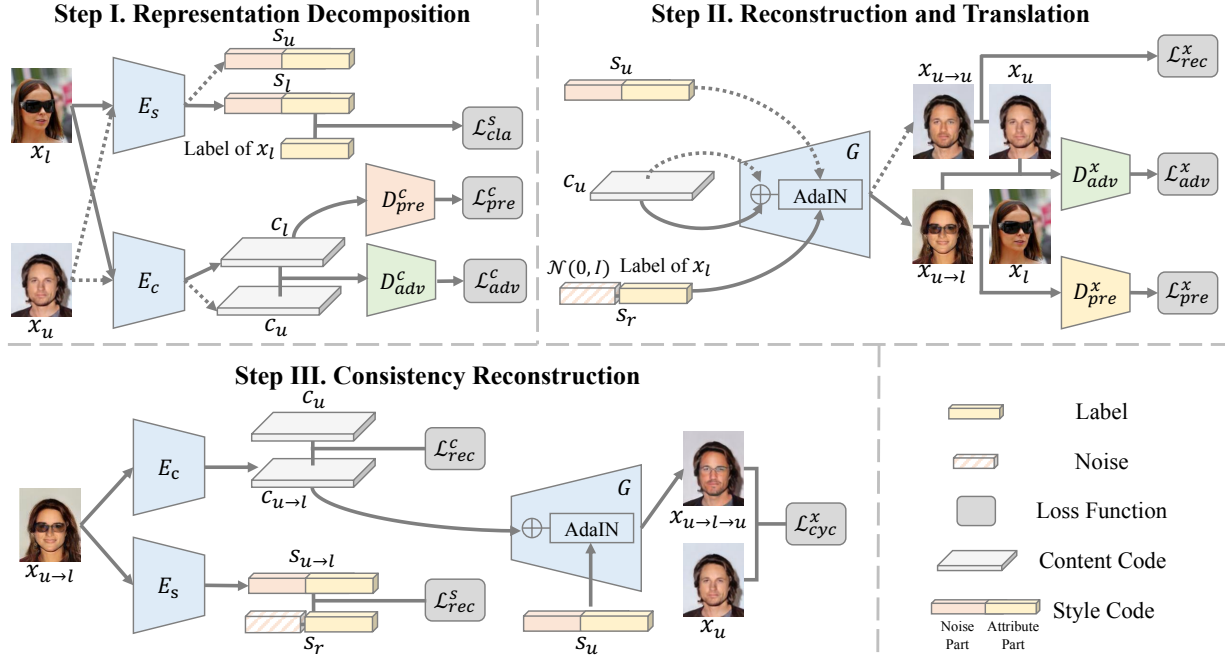
**Step I. Representation Decomposition**

**Step II. Reconstruction and Translation**

**Step III. Consistency Reconstruction**

Figure 2. The flowchart of training AGUIT. First, the representation of image $x_l$ and image $x_u$ is decomposed into style code $s_l, s_u$ and content code $c_l, c_u$ by style encoder $E_s$ and content encoder $E_c$. The label $l = [1, -1, 1, ...]$ corresponds to the attributes [glasses, male, smile,...] and 1 means $x_l$ has the corresponding attribute, -1 is opposite. Second, the generator $G$ decodes $s_u, c_u$ to reconstruct $x_u$, and decodes $s_r, c_u$ to translate attributes of $x_l$ to $x_u$. The $s_r$ consists of noise $z \sim \mathcal{N}(0, I)$ and label $l$. The reconstructed and translated images are denoted as $x_{u \to u}$ and $x_{u \to l}$, respectively. Third, $x_u$ is cyclically reconstructed by using $c_{u \to l}$ and $s_u$.

models. The works of [35, 31] train a discriminator to classify input into different classes. The work of [9] introduces a separate discriminator and classifier models. Other approaches incorporate inference models to predict missing labeled features [10] or harness the joint distribution of labels and data matching [11]. Recently, the setting of training a classifier from a few labels is introducing to generative model. The work of [28] learns a generative model from a few labels by clustering deep features. Unlike the above works, we focus on exploiting semi-supervised learning for image-to-image translation task.

**Image-to-Image Translation.** Image-to-image translation tasks have attracted increasing focus. For instance, Pix2pix [17] achieves translation based on paired image. To use unpaired images, CycleGAN [40], DiscoGAN [19], DualGAN [37] and UNIT [24] are proposed. For multi-domain translation, ComboGAN [2] learns multiple generators and StarGAN [8] reduces them to a single one by inputting the target labels and images together. GANimation [32] describes the labels in a continuous manifold. CerfGAN [25] adopts a multi-class discriminator to enable the generator to translate images with high domain shifts. For multi-modal translation, BicycleGAN [41] extends pix2pix by learning a stochastic mapping from source to target. MUNIT [16] and DRIT [22] decompose the image representation into style code and content code, and then decode back to translated

images. Augmented CycleGAN [1] learns many-to-many translation by using the stochastic mappings. Recently, SMIT [33] is proposed to solve these basic tasks jointly, in which the images with labels and additional style noises are input to the generator. Unlike existing methods, AGUIT encodes the attributes into the style code to conduct high-level translation tasks.

**Representation Disentanglement.** Representation disentanglement aims at disentangling and explaining the learned representation. The methods can be divided into two categories, *i.e.*, supervised disentanglement and unsupervised disentanglement. Supervised disentanglement methods [7, 30, 29, 20, 26, 23] make use of labeled data while unsupervised disentanglement methods [6, 13, 3, 18, 5] learn the properties from unlabeled data. The proposed AGUIT can conduct both supervised and unsupervised disentanglement by the proposed ingenious style code.

## 3. The Model of AGUIT

### 3.1. Problem Formulation

Let $(\mathcal{X}_l, L)$ denote the pairs of images with corresponding attribute labels and $\mathcal{X}_u$ denote the images without attribute labels. The image representation is decomposed into the style code $\mathcal{S}$ and the content code $\mathcal{C}$. Our goal is to train a model for UIT by using $(\mathcal{X}_l, L)$ and $\mathcal{X}_u$, in which $\mathcal{C}$ retains

the content of objects and $\mathcal{S}$ learns disentanglement for attributes. After training, by manipulating $\mathcal{S}$, the expected attributes should be translated to the outputs. The flowchart of training the proposed AGUIT for tackling this problem is shown in Fig. 2.

## 3.2. Representation Decomposition

As the step I in Fig. 2, we decompose the image representation of $\mathcal{X}_l$ into $\mathcal{S}_l$ and $\mathcal{C}_l$ by style encoder $E_s$ and content encoder $E_c$. The same operation is conducted on $\mathcal{X}_u$ to get $\mathcal{S}_u$ and $\mathcal{C}_u$. Notably, $\mathcal{S}$ consists of the noise part and the attribute part.

**Style Classifying Loss.** To encourage continuity for the attribute part of style code, we enforce it to be close to $L$ in the Euclidean space. This design enhances the capability of supervised disentanglement, because the continuous code can be used for fine control. The style classifying loss is defined as:

$$\mathcal{L}_{cla}^s = \mathbb{E}_{(x_l, l) \in (\mathcal{X}_l, L)}[\|E_s(x_l)_l - l\|_2], \qquad (1)$$

where $E_s(x_l)_l$ denotes the attribute part of style code.

**Content Confusing Loss.** To constrain the content code to a common space, we utilize a discriminator $D_{adv}^c$ to train $E_c$. In the training process, $D_{adv}^c$ tries to distinguish $\mathcal{C}_l$ and $\mathcal{C}_u$, while $E_c$ learns a common representation for them. The content confusing loss is as follows:

$$\mathcal{L}_{adv}^c = \mathbb{E}_{x_l \in \mathcal{X}_l, x_u \in \mathcal{X}_u}[\log(D_{adv}^c(E_c(x_u))) \\ + \log(1 - D_{adv}^c(E_c(x_l)))]. \qquad (2)$$

**Content-Style Separating Loss.** To make the style code and content code independent to each other, we introduce a predictor $D_{pre}^c$ as inspired by PM [34]. Because the style code is not stable in the early training phase, we let $D_{pre}^c$ directly predict the image labels conditioned on content code and prevent $E_c$ from learning the information of image labels. Therefore, the content-style separating loss is defined as:

$$\mathcal{L}_{pre}^c = \mathbb{E}_{(x_l, l) \in (\mathcal{X}_l, L)}[\log(D_{pre}^c(l|E_c(x_l)))]. \qquad (3)$$

## 3.3. Reconstruction and Translation

As the step II in Fig. 2, we reconstruct $\mathcal{X}_u$ by inputting $\mathcal{S}_u$ and $\mathcal{C}_u$ to the generator $G$ with AdaIN [15]. Similarly, the translation of $\mathcal{X}_u$ is done by inputting $\mathcal{S}_r$ and $\mathcal{C}_u$ to $G$. The $\mathcal{S}_r$ consists of the random noise $Z \sim \mathcal{N}(0, I)$ and the attribute label $L$ of $\mathcal{X}_l$. We denote the reconstructed images as $\mathcal{X}_{u \rightarrow u}$ and the translated images as $\mathcal{X}_{u \rightarrow l}$.

**Image Reconstructing Loss.** To guarantee that the learned representation can be decoded to generate target images, we use an image reconstructing loss to train the translator (*i.e.*, $E_s$, $E_c$ and $G$):

$$\mathcal{L}_{rec}^x = \mathbb{E}_{x_u \in \mathcal{X}_u}[\|x_u - G(E_c(x_u), E_s(x_u))\|_1]. \qquad (4)$$

**Algorithm 1** Process of Training AGUIT

**Input:** The images with domain labels: $(\mathcal{X}_l, L)$. The images without requirement of labels: $\mathcal{X}_u$.

**Output:** The learned style encoder $E_s$, content encoder $E_c$, and generator $G$.

1: **while** not convergence **do**
2:     Get $\mathcal{S}_l, \mathcal{C}_l$ by $E_s(\mathcal{X}_l), E_c(\mathcal{X}_l)$.
3:     Get $\mathcal{S}_u, \mathcal{C}_u$ by $E_s(\mathcal{X}_u), E_c(\mathcal{X}_u)$.
4:     Reconstruct $\mathcal{X}_u$: $\mathcal{X}_{u \rightarrow u} = G(\mathcal{C}_u, \mathcal{S}_u)$.
5:     Construct random style code: $\mathcal{S}_r = [\mathcal{N}(0, I), L]$.
6:     Translate $\mathcal{X}_u$ under style $\mathcal{S}_r$: $\mathcal{X}_{u \rightarrow l} = G(\mathcal{C}_u, \mathcal{S}_r)$.
7:     Get $\mathcal{S}_{u \rightarrow l}, \mathcal{C}_{u \rightarrow l}$ by $E_s(\mathcal{X}_{u \rightarrow l}), E_c(\mathcal{X}_{u \rightarrow l})$.
8:     Cycle reconstruct $\mathcal{X}_u$: $\mathcal{X}_{u \rightarrow l \rightarrow u} = G(\mathcal{C}_{u \rightarrow l}, \mathcal{S}_u)$.
9:     Let $A$ denote $E_s, E_c, G$.
10:    Let $B$ denote $D_{adv}^c, D_{adv}^x, D_{pre}^c, D_{pre}^x$.
11:    Fixing $A$, optimize $B$ by Eq. 11.
12:    Fixing $B$, optimize $A$ by Eq. 10.
13: **end while**
14: **return** $E_s, E_c$ and $G$.

**Image Adversarial Loss.** To make the translator translate attributes $L$ to $\mathcal{X}_u$ while retaining the content, we use a discriminator $D_{adv}^x$ which attempts to distinguish $\mathcal{X}_u$ and $\mathcal{X}_{u \rightarrow l}$. The translator tries to fool $D_{adv}^x$ by generating $\mathcal{X}_{u \rightarrow l}$ with the same content of $\mathcal{X}_u$. Then, the image adversarial loss is defined as:

$$\mathcal{L}_{adv}^x = \mathbb{E}_{x_u \in \mathcal{X}_u, x_l \in \mathcal{X}_l}[\log(D_{adv}^x(G(E_c(x_u), s_r))) \\ + \log(1 - D_{adv}^x(x_l))], \qquad (5)$$

where $s_r \in \mathcal{S}_r$ contains a random noise $z \sim \mathcal{N}(0, I)$ and the attribute label $l \in L$ of $x_l$.

**Image Classifying Loss.** To make $\mathcal{X}_{u \rightarrow l}$ have the same attributes as $\mathcal{X}_l$, we apply a classifier $D_{pre}^x$ inspired by AC-GAN [31]. In the process, $D_{pre}^x$ attempts to predict $\mathcal{X}_l$ under the supervision of $L$, and the translator also tries to generate images satisfying the given label $L$. Therefore, the domain classifying loss for $D_{pre}^x$ is defined as:

$$\mathcal{L}_{pre}^{x,D} = \mathbb{E}_{(x_l, l) \in (\mathcal{X}_l, L)}[\log(D_{pre}^x(l|x_l))]. \qquad (6)$$

The loss for translator is defined as:

$$\mathcal{L}_{pre}^{x,G} = \mathbb{E}_{x_u \in \mathcal{X}_u}[\log(D_{pre}^x(l|G(E_c(x_u), s_r)))]. \qquad (7)$$

## 3.4. Consistency Reconstruction

As the step III in Fig. 2, we cyclically reconstruct the unlabeled images $\mathcal{X}_u$ with the style code $\mathcal{S}_{u \rightarrow l}$ and content code $\mathcal{C}_{u \rightarrow l}$ of $\mathcal{X}_{u \rightarrow l}$ encoded by $E_s$ and $E_c$.

**Cycle Consistency Loss.** To guarantee that $\mathcal{X}_{u \rightarrow l}$ preserves the content of $\mathcal{X}_u$ while changing only the domain-specific style of label $L$, we apply a cycle consistency loss

to the translator, which is defined as:

$$\mathcal{L}_{cyc}^x = \mathbb{E}_{x_u \in \mathcal{X}_u}[\|x_u - G(E_c(G(E_c(x_u), s_r)),$$
$$E_s(x_u))\|_1]. \quad (8)$$

**Feature Consistency Loss.** To preserve the consistency of $\mathcal{C}_u$ and $\mathcal{S}_u$ in the representation level, we apply a feature consistency loss [41] to the translator as follows:

$$\mathcal{L}_{lat} = \mathcal{L}_{rec}^c + \mathcal{L}_{rec}^s$$
$$= \mathbb{E}_{x_u \in \mathcal{X}_u}[\|E_c(x_u) - E_c(G(E_c(x_u), s_r))\|_1 \quad (9)$$
$$+ \|s_r - E_s(G(E_c(x_u), s_r))\|_1]].$$

### 3.5. Optimization and Inference for AGUIT

In the training phase, we optimize the style encoder $E_s$, content encoder $E_c$, generator $G$, predictors $D_{pre}^c$, $D_{pre}^x$ and discriminators $D_{adv}^c$, $D_{adv}^x$ for AGUIT jointly. We can write the overall objective of $E_s$, $E_c$ and $G$ as:

$$\mathcal{L}_{G,E_c,E_s} = \lambda_{cla}^s \mathcal{L}_{cla}^s + \lambda_{adv}^c \mathcal{L}_{adv}^c - \lambda_{pre}^c \mathcal{L}_{pre}^c$$
$$+ \lambda_{rec}^x \mathcal{L}_{rec}^x + \lambda_{adv}^x \mathcal{L}_{adv}^x + \lambda_{pre}^{x,G} \mathcal{L}_{pre}^{x,G} \quad (10)$$
$$+ \lambda_{cyc}^x \mathcal{L}_{cyc}^x + \lambda_{lat} \mathcal{L}_{lat}.$$

The objective of $D_{adv}^c$, $D_{adv}^x$ and $D_{pre}^c$, $D_{pre}^x$ is defined as:

$$\mathcal{L}_D = -\lambda_{adv}^c \mathcal{L}_{adv}^c + \lambda_{pre}^c \mathcal{L}_{pre}^c$$
$$- \lambda_{adv}^x \mathcal{L}_{adv}^x + \lambda_{pre}^{x,D} \mathcal{L}_{pre}^{x,D}, \quad (11)$$

where $\lambda$ with different superscript and subscript are hyperparameters for balancing the losses. Alg. 1 summarizes the training process of AGUIT.

In the inference phase, a test image $x_t$ is encoded into $s_t$ and $c_t$ by $E_s$ and $E_c$. Then, the attributes of $x_t$ can be changed by manipulating one or more dimensions of the style code $s_t$ optionally. The $c_t$ and the manipulated style code $s_t'$ are input into the generator $G$ and the manipulated style is translated while the content of original input $x_t$ is retained[1].

## 4. Experiments

### 4.1. Experimental Settings

**Datasets.** We evaluate AGUIT on the Dog2Cat [22] dataset which contains a variety of cat and dog faces, as well as the CelebA [27] dataset which contains more than 200,000 labeled faces with attributes such as hair color, gender and presence of eyeglasses. Among 40 attributes of CelebA dataset, we choose 8 most common attributes (*i.e.*, Black Hair, Blond Hair, Brown Hair, Gender, Age, Smiling, Eyeglasses, Goatee) for experiments.

---

[1]More detailed description for the inference of AGUIT can be found in our supplementary materials.
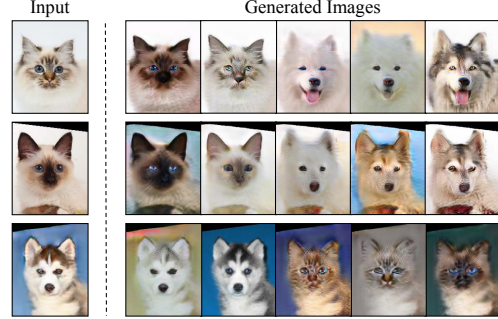


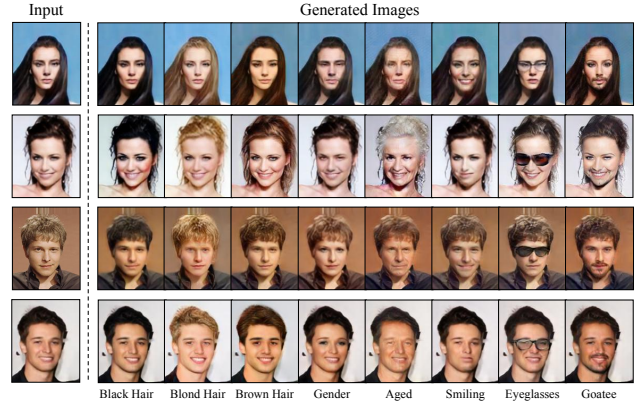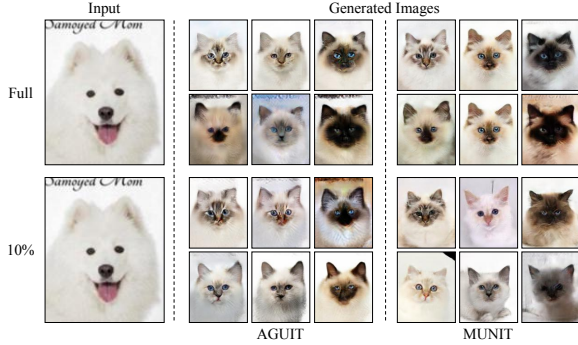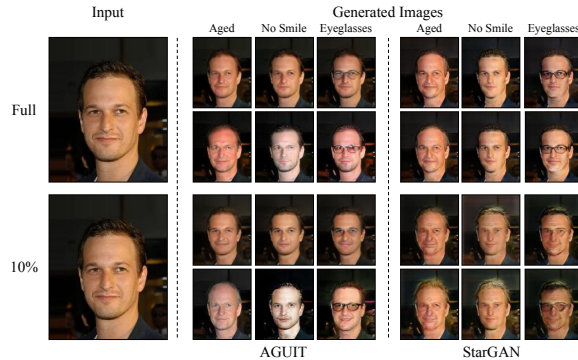Figure 3. Examples of multi-modal translation with AGUIT on Dog2Cat dataset.



Figure 4. Examples of multi-domain translation with AGUIT on CelebA dataset.

**Evaluation Metrics.** Firstly, to compare the quality of translated images, we conducted *human preference* on Amazon Mechanical Turk (AMT) for evaluation. The workers were given a source image with two target images translated by different models and required to answer the following questions, *i.e.*, which target image has more precise attributes of the offered label, which target image retains more similar shape of source image and which target image looks more natural. We randomly chose 100 source images. For multi-modal translation, each source image has 19 random translations. For multi-label translation, each source image has 30 label-specific translations. Secondly, to measure the diversity, we compute the average *LPIPS Distance* [39] between pairs of randomly sampled translation outputs with the same settings as [41].

**Baselines.** CycleGAN [40] consists of two residual translation networks trained with adversarial loss and cycle reconstruction loss. MUNIT [16] or DRIT [22] decompose the image representation into content code and style code. Then the translated image is generated by recombining the content code with a random style code sampled from the style space of target domain. StarGAN [8] uses the images with labels as the input and manipulates the label to translate the corresponding attributes on the output images.

(a) Multi-modal translation



(b) Multi-domain translation.

Figure 5. Examples of basic translation tasks under semi-unsupervised setting.

**Implementation Details.** We use the input image size of $128 \times 128$ for all our experiments. More results of $128 \times 128$ and $512 \times 512$ inputs, detailed settings of hyper-parameters, network architectures and optimizers can be found in our supplementary materials.

### 4.2. Basic Image Translation Tasks

We show that AGUIT can do basic image translation tasks (*i.e.*, multi-modal translation and multi-domain translation) under the unpaired setting.

For multi-modal translation, the results are shown in Fig. 3, from which we can see that AGUIT produces a variety of results with the same or different species. The translated images greatly retain the content (*i.e.*, shape and posture) of input images.

For multi-domain translation, we change one of the labeled attributes, and the outputs are shown in Fig. 4. From the results we can see that AGUIT learns attributes for the style code under the supervision of domain labels.

### 4.3. Benefits of Semi-supervised Setting

In this section, we show the effectiveness of AGUIT under the semi-supervised setting, in which the training set contains labeled images mixed with unlabeled images. We compare AGUIT with MUNIT and StarGAN under the set-

| (a) Multi-modal translation | | | (b) Multi-domain translation | | |
|---|---|---|---|---|---|
| Method-Settings | Quality | Diversity | Method-Settings | Quality | Diversity |
| MUNIT[16]-10%OL | 9.3% | **0.519** | StarGAN[8]-10%OL | 20.2% | 0.186 |
| AGUIT-10%OL | 13.3% | 0.327 | AGUIT-10%OL | 23.8% | 0.231 |
| AGUIT-10% | **27.0%** | 0.442 | AGUIT-10% | **30.8%** | **0.252** |
| MUNIT[16]-full | 24.9% | **0.492** | StarGAN[8]-full | 34.7% | 0.165 |
| AGUIT-full | **50.0%** | 0.438 | AGUIT-full | **50.0%** | **0.242** |

Table 2. Quantitative evaluations of basic translation tasks under semi-supervised setting. The column of 'Diversity' is the average LPIPS distance. The column of 'Quality' is the human preference score. The setting of 10%OL means only 10% of training images are used. The setting of 10% means 90% training images are treated as unlabeled data for training.
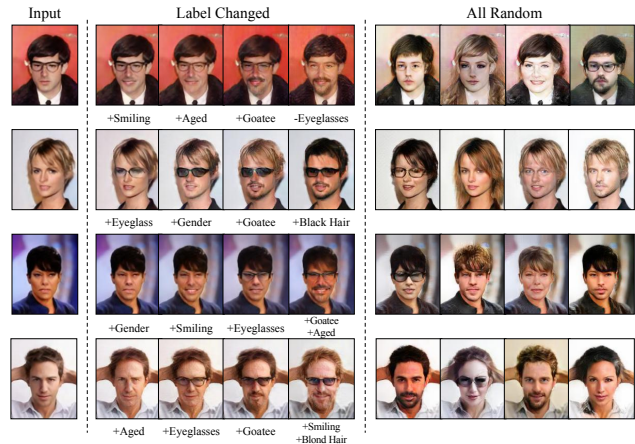


Figure 6. Examples of style specific translation. The content codes of input images remain unchanged. From the second column to the fifth column, we gradually add changes to the attribute part of style code. From the sixth column to the ninth column, we randomly change all dimensions of the style code.

tings of 10% labels and full labels for multi-modal translation and multi-domain translation. For AGUIT, we randomly choose images from labeled images to form the training pairs under the fully label setting.

**Qualitative Evaluation.** As shown in Fig. 5(a), for multi-modal translation, AGUIT and MUNIT both have good results when using fully labeled images for training. But when the number of labeled training images decreases drastically, the content of images generated by MUNIT change a lot compared to the input image. On the contrary, AGUIT still generates a variety of results retaining the shape of input image. As shown in Fig. 5(b), for multi-domain translation, AGUIT and StarGAN both translate attributes well when using fully labeled images as the training set. However, StarGAN cannot generate diversified results. When the number of labeled training images decreases drastically, StarGAN fails both on the quality and diversity. On the contrary, AGUIT accomplishes both goals well, as unlabeled images can be incorporated into the training phase and the data distribution is augmented.

**Style Limited Translation**

Input | Limitation | Generated Images

**Style Manipulative Translation**

Input | Reference | Manipulation | Generated Images

keep gender.

Keep age and hair color.

Don't smile

Keep all domain label.

All style random.

Just gender.

Transfer age and hair color.

Don't change the smiling but all domain label.

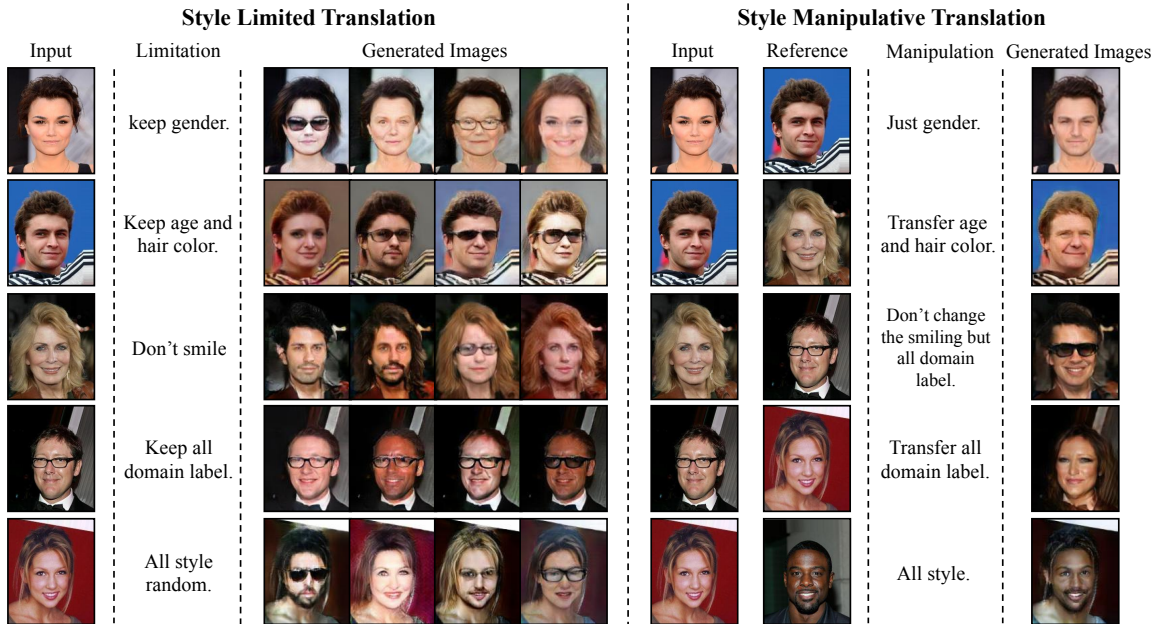Transfer all domain label.

All style.

Figure 7. Examples of style limited translation (left) and style manipulative translation (right). On the style limited translation, the content and style codes are given by the input images. Some limitations are applied to the style code, and various output images are produced by AGUIT. On the style manipulative translation, the content and style codes are given by the input images and style reference images, respectively. Then manipulations are applied to the style code to translate the input images with the manipulated style.

**Quantitative Evaluation.** As shown in Tab. 2(a), for multi-modal translation, the human preference score of AGUIT is far beyond MUNIT, although MUNIT has a high value on LPIPS distance. Commonly, high human preference score means high quality, and high LPIPS distance means high diversity. We set high quality as prior to high diversity due to the goal of retaining the shape of input images. From this perspective, AGUIT performs well compared with MUNIT. As shown in Tab. 2(b), for multi-domain translation, AGUIT is superior to StarGAN on the evaluations of both quality and diversity. For both tasks, the performance of AGUIT increases after incorporating unlabeled data. To explain, the extra information of unlabeled data enriches the learning of representation.

Overall, based on the above qualitative and quantitative evaluations, AGUIT has stronger capability on the basic image translation tasks, and gains good benefits from the semi-supervised scheme.

### 4.4. Style Operable Image Translation Tasks

Because the style information is encoded into the style code that contains two variables drawn from the standard Gaussian distribution (*i.e.*, the noise part) and the attributes distribution (*i.e.*, the attribute part), AGUIT can accomplish high-level tasks and produce more plentiful results.

**Style Specific Translation.** As shown in Fig. 6, we translate input image via style code, whose attribute part is

assigned by changing specific dimensions. The labeled attributes can be assigned cumulatively or randomly for translating input images. The results show AGUIT translates images accurately with different specified style codes, while retaining the content of the input images.

**Style Limited Translation.** As shown in Fig. 7 left, we translate images via giving some limitations to the style code of input image. From the results we can see AGUIT produces various results to meet the limitations. In fact, an operation (*e.g.* reverse, hold, random sample) can be applied to any dimensions of style codes unrelated with the limitations to enrich the output images.

**Style Manipulative Translation.** As shown in Fig. 7 right, we translate input images under style reference images with some manipulations. The style reference image provides style code. Then the manipulations can be assigned to it for translation.

Overall, the above three experiments show that AGUIT has the ability to handle the style operable image translation tasks and has great flexibility.

### 4.5. Fine-controlled Image Translation Tasks

In order to further reveal the ability of AGUIT, we conduct interpolation on both attribute part and noise part of style code for evaluating the effectiveness of disentanglement for AGUIT. The results show that the representation learned by AGUIT is continuous and disentangled on both
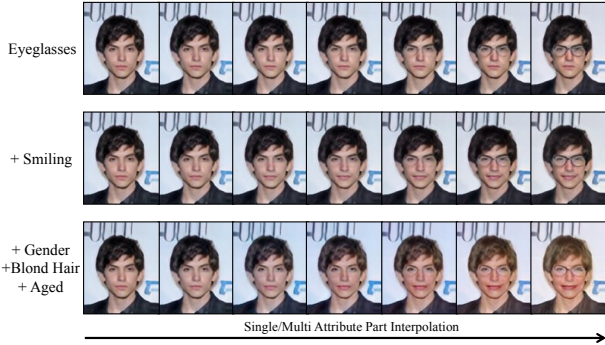
Figure 8. Supervised disentanglement on the attribute part of style code. The interpolations are conducted from left to right. Only one attribute is operated on the first row, and more attributes are added for the following rows.
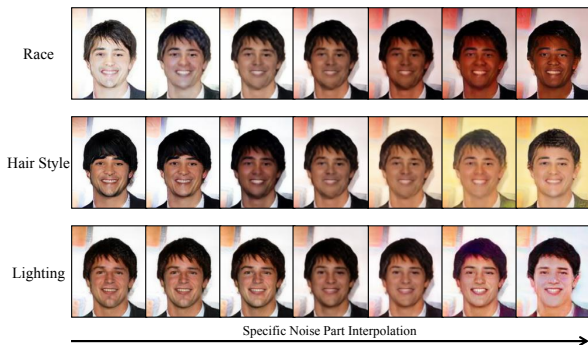


Figure 9. Unsupervised disentanglement on the noise part of style code. The interpolations are conducted on single dimension from left to right on every row.

supervised attribute part and unsupervised noise part. The outputs can be fine-controlled by adjusting the style code gradually.

**Supervised Disentanglement.** As shown in Fig. 8, we can see that the labeled attributes are disentangled and can be fine-controlled while interpolating the value of the corresponding dimension. The results suggest that whatever the single or multiple interpolation on the attribute part of style code, the representation space is independent and continuous. The clean translated images indicate that the supervised disentanglement of style code has very high quality.

**Unsupervised Disentanglement.** Similarly, as shown in Fig. 9, we can see that some unlabeled attributes such as lighting and hair style are learned and disentangled by AGUIT. Though the background contains a little bit perturbation, the unsupervised disentanglement, to some extent, is a possible way to reduce the requirement of the attributes need to be labeled.

Overall, from the above experiments, AGUIT can achieve fine control for outputs because of the disentangled style code.
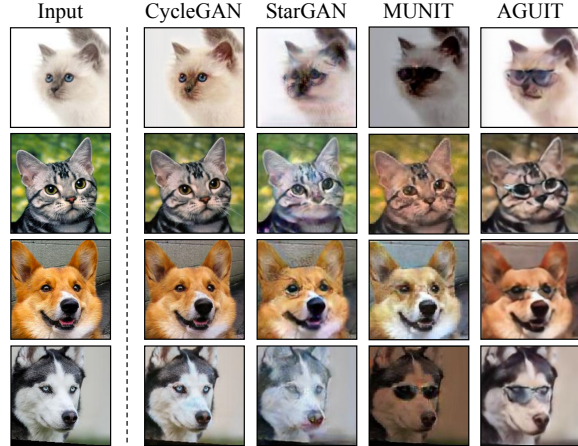


Figure 10. Example results of disentangled transfer.

## 4.6. Disentangled Transfer

We further introduce a new task for image-to-timage translation, termed disentangled transfer. Disentanglement is one of the goals of representation learning. It is of great benefit for the practical application if the information learned by disentangled representation can be transferred. For example, we can create zero-shot samples which has unrelated attributes.

To illustrate the goal of disentangled transfer clearly, we give an example as follows: We utilize models trained on CelebA dataset to make the cats and dogs wear eyeglasses. Because the eyes of dogs and cats are similar to human's, the attribute for human eyes ought to be transferred if the representation is disentangled well. For comparison, we train CycleGAN and MUNIT by human faces with or without eyeglasses in CelebA dataset, and train the StarGAN by full CelebA dataset. As shown in Fig. 10, we can see that CycleGAN cannot put the eyeglasses to the face of dogs and cats. We attribute it to the poor representation learning of CycleGAN. StarGAN also fails on this task and brings noise to the translated images. We attribute it to the fact that StarGAN is not capable for representation learning of attributes. MUNIT has better results than the above models, where some translations achieve the goal while some are not, but the hue of the successful examples is changed drastically. We attribute it to that the representation is not disentangled well, so that small turbulence on one dimension can change the whole image via the decoding scheme.

On the contrary, AGUIT works very well in this task. Although it has a little changes on background, AGUIT finds the position of the eyes in the input image and puts the glasses on. It suggests that AGUIT learns a good representation than other models, which also indicates the good potential of the disentangled transfer for image-to-image translation.

8

## 5. Conclusion

In this work, we propose a model termed as AGUIT, which is the first model to consider UIT tasks with a semi-supervised setting, and achieves disentanglement of representation as well as fine control of outputs for UIT. The semi-supervised learning scheme and the encoding of style information bring the capability to AGUIT. Extensive experiments demonstrate the performance of AGUIT over the state-of-the-art image-to-image translation models. First, we showed that AGUIT can do the basic image translation tasks (*i.e.*, multi-modal and multi-domain translations). Second, we qualitatively and quantitatively evaluated the benefits brought by the semi-supervised setting of AGUIT. Third, we exhibited the style manipulation and fine control results. Finally, we revealed that AGUIT can carry out disentangled transfer, which is a new task introduced to image-to-image translation.

## References

[1] A. Almahairi, S. Rajeswar, A. Sordoni, P. Bachman, and A. Courville. Augmented cyclegan: Learning many-to-many mappings from unpaired data. In *ICML*, 2018. 3

[2] A. Anoosheh, E. Agustsson, R. Timofte, and L. Van Gool. Combogan: Unrestrained scalability for image domain translation. In *CVPR*, 2018. 3

[3] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, and A. Lerchner. Understanding disentangling in beta-vae. In *NIPS*, 2018. 3

[4] Q. Chen and V. Koltun. Photographic image synthesis with cascaded refinement networks. In *ICCV*, 2017. 1

[5] T. Q. Chen, X. Li, R. B. Grosse, and D. K. Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *NIPS*, 2018. 3

[6] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *NIPS*, 2016. 3

[7] B. Cheung, J. A. Livezey, A. K. Bansal, and B. A. Olshausen. Discovering hidden factors of variation in deep networks. *arXiv preprint arXiv:1412.6583*, 2014. 3

[8] Y. Choi, M. Choi, M. Kim, J. W. Ha, S. Kim, and J. Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, 2018. 1, 2, 3, 5, 6

[9] L. Chongxuan, T. Xu, J. Zhu, and B. Zhang. Triple generative adversarial nets. In *NIPS*, 2017. 3

[10] E. L. Denton et al. Unsupervised learning of disentangled representations from video. In *NIPS*, 2017. 3

[11] Z. Gan, L. Chen, W. Wang, Y. Pu, Y. Zhang, H. Liu, C. Li, and L. Carin. Triangle generative adversarial networks. In *NIPS*, 2017. 3

[12] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014. 2

[13] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017. 3

[14] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. A. Efros, and T. Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *ICML*, 2018. 1

[15] X. Huang and S. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017. 4

[16] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz. Multimodal unsupervised image-to-image translation. In *ECCV*, 2018. 2, 3, 5, 6

[17] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 1, 2, 3

[18] H. Kim and A. Mnih. Disentangling by factorising. *arXiv preprint arXiv:1802.05983*, 2018. 3

[19] T. Kim, M. Cha, H. Kim, J. K. Lee, and J. Kim. Learning to discover cross-domain relations with generative adversarial networks. In *ICML*, 2017. 3

[20] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling. Semi-supervised learning with deep generative models. In *NIPS*, 2014. 3

[21] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017. 1

[22] H.-Y. Lee, H.-Y. Tseng, J.-B. Huang, M. Singh, and M.-H. Yang. Diverse image-to-image translation via disentangled representations. In *ECCV*, 2018. 2, 3, 5

[23] A. H. Liu, Y.-C. Liu, Y.-Y. Yeh, and Y.-C. F. Wang. A unified feature disentangler for multi-domain image translation and manipulation. In *NIPS*, 2018. 3

[24] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *NIPS*, 2017. 1, 3

[25] X. Liu, S. Zhang, H. Liu, X. Liu, and R. Ji. Cerfgan: A compact, effective, robust, and fast model for unsupervised multi-domain image-to-image translation. *arXiv preprint arXiv:1805.10871v2*, 2018. 3

[26] Y. C. Liu, Y. Y. Yeh, T. C. Fu, W. C. Chiu, and Y. C. F. Wang. Detach and adapt: Learning cross-domain disentangled deep representation. In *CVPR*, 2018. 3

[27] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *ICCV*, 2015. 5

[28] M. Lucic, M. Tschannen, M. Ritter, X. Zhai, O. Bachem, and S. Gelly. High-fidelity image generation with fewer labels. *arXiv preprint arXiv:1903.02271*, 2019. 3

[29] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015. 3

[30] M. F. Mathieu, J. J. Zhao, J. Zhao, A. Ramesh, P. Sprechmann, and Y. LeCun. Disentangling factors of variation in deep representation using adversarial training. In *NIPS*, 2016. 3

[31] A. Odena, C. Olah, and J. Shlens. Conditional image synthesis with auxiliary classifier gans. In *ICML*, 2017. 3, 4

[32] A. Pumarola, A. Agudo, A. M. Martinez, A. Sanfeliu, and F. Moreno-Noguer. Ganimation: Anatomically-aware facial animation from a single image. In *ECCV*, 2018. 2, 3

[33] A. Romero, P. Arbeláez, L. Van Gool, and R. Timofte. Smit: Stochastic multi-label image-to-image translation. *arXiv preprint arXiv:1812.03704*, 2018. 2, 3

[34] J. Schmidhuber. Learning factorial codes by predictability minimization. *Neural Computation*, 1992. 4

[35] J. T. Springenberg. Unsupervised and semi-supervised learning with categorical generative adversarial networks. *arXiv preprint arXiv:1511.06390*, 2015. 3

[36] B. Wu, H. Duan, Z. Liu, and G. Sun. Srpgan: Perceptual generative adversarial network for single image super resolution. In *CVPR*, 2017. 1

[37] Z. Yi, H. Zhang, P. Tan, and M. Gong. Dualgan: Unsupervised dual learning for image-to-image translation. In *ICCV*, 2017. 1, 3

[38] R. Zhang, P. Isola, and A. A. Efros. Colorful image colorization. In *ECCV*, 2016. 1

[39] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5

[40] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 1, 2, 3, 5

[41] J.-Y. Zhu, R. Zhang, D. Pathak, T. Darrell, A. A. Efros, O. Wang, and E. Shechtman. Toward multimodal image-to-image translation. In *NIPS*, 2017. 2, 3, 5

# A. Implement Details

**Architecture.** The network architectures of AGUIT are shown as follows in Fig. 11. K: kernel size, S: stride size, P: padding size, IN: instance normalization. LN: layer normalization, AdaIN: adaptive instance normalization, nd: the dimension of attribute part, nz: the dimension of noise part.
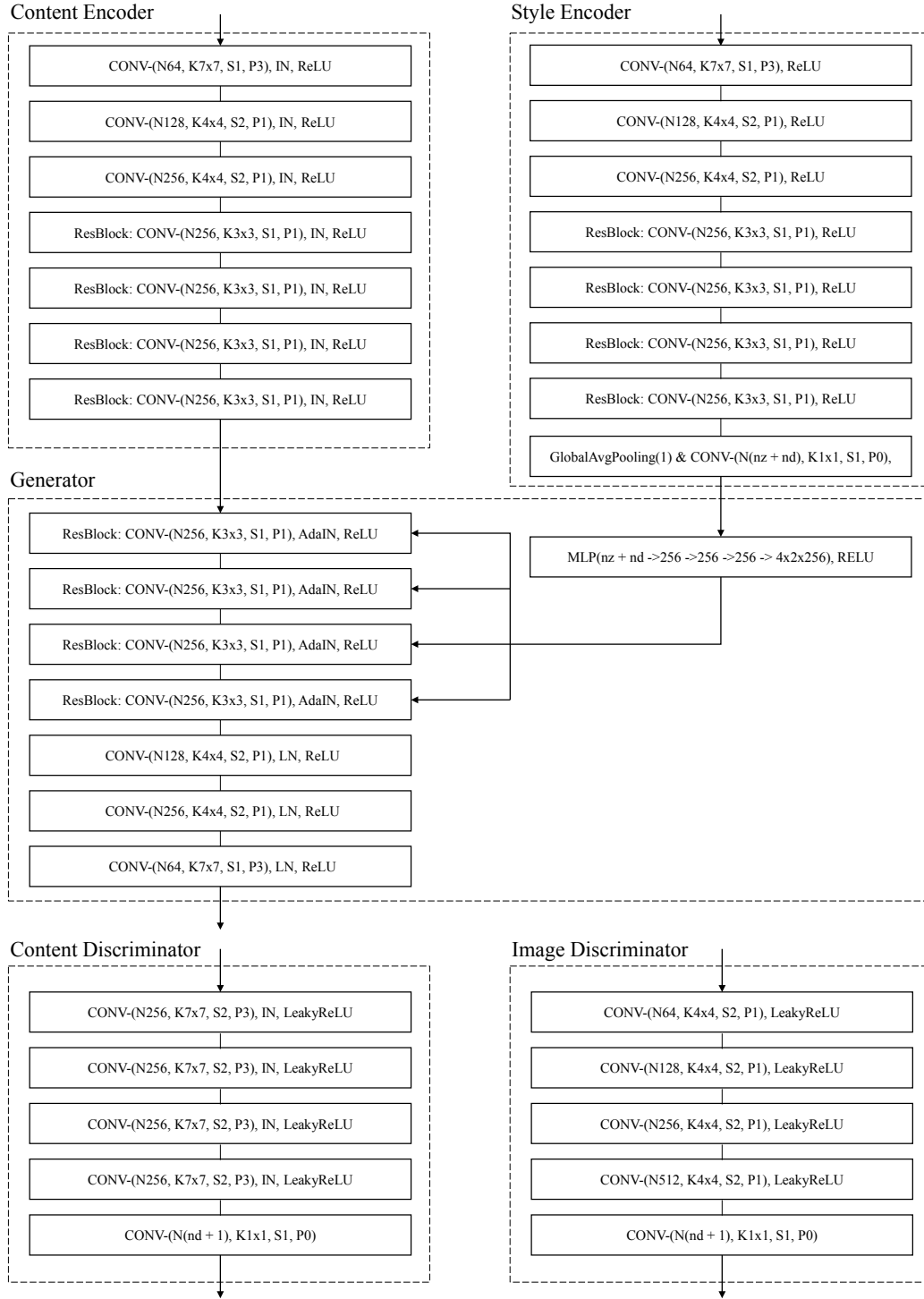
Content Encoder

| CONV-(N64, K7x7, S1, P3), IN, ReLU |
| CONV-(N128, K4x4, S2, P1), IN, ReLU |
| CONV-(N256, K4x4, S2, P1), IN, ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), IN, ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), IN, ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), IN, ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), IN, ReLU |

Style Encoder

| CONV-(N64, K7x7, S1, P3), ReLU |
| CONV-(N128, K4x4, S2, P1), ReLU |
| CONV-(N256, K4x4, S2, P1), ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), ReLU |
| GlobalAvgPooling(1) & CONV-(N(nz + nd), K1x1, S1, P0), |

Generator

| ResBlock: CONV-(N256, K3x3, S1, P1), AdaIN, ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), AdaIN, ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), AdaIN, ReLU |
| ResBlock: CONV-(N256, K3x3, S1, P1), AdaIN, ReLU |
| CONV-(N128, K4x4, S2, P1), LN, ReLU |
| CONV-(N256, K4x4, S2, P1), LN, ReLU |
| CONV-(N64, K7x7, S1, P3), LN, ReLU |

| MLP(nz + nd ->256 ->256 ->256 -> 4x2x256), RELU |

Content Discriminator

| CONV-(N256, K7x7, S2, P3), IN, LeakyReLU |
| CONV-(N256, K7x7, S2, P3), IN, LeakyReLU |
| CONV-(N256, K7x7, S2, P3), IN, LeakyReLU |
| CONV-(N256, K7x7, S2, P3), IN, LeakyReLU |
| CONV-(N(nd + 1), K1x1, S1, P0) |

Image Discriminator

| CONV-(N64, K4x4, S2, P1), LeakyReLU |
| CONV-(N128, K4x4, S2, P1), LeakyReLU |
| CONV-(N256, K4x4, S2, P1), LeakyReLU |
| CONV-(N512, K4x4, S2, P1), LeakyReLU |
| CONV-(N(nd + 1), K1x1, S1, P0) |

Figure 11. Detailed architecture of AGUIT.

**Training Details.** For training, we use the Adam optimizer with a batch size of 8 for 128x128 photos of CelebA dataset and 1 for both Dog2Cat dataset and 512x512 photos of CelebA dataset. The learning rate is 0.0001, and exponential decay rates is set as $(\beta_1, \beta_2) = (0.5, 0.999)$. In all experiments, we set the hyper-parameters as follows: $\lambda^s_{cla} = 10$, $\lambda^c_{adv} = 1$, $\lambda^c_{pre} = 1$, $\lambda^x_{rec} = 10$, $\lambda^x_{adv} = 1$, $\lambda^x_{pre} = 1$, $\lambda^x_{cyc} = 10$ and $\lambda_{lat} = 10$, we use LSGAN for all GAN loss.

## B. Additional Experiment

The inference phase of AGUIT is shown in Fig. 12. The operations for style code includes: Hold: the same as original, Reverse: the opposite to original, Random: the random change, Replace: the replacement to a reference style code of another image, Value: the value to a specific number.

In Fig. 13, we show some additional results for style operable image translation tasks in CelebA dataset. In Fig. 14 and Fig. 15, we show some additional results for fine-controlled image translation tasks with both noise part and attribute part of style code. 512x512 results are shown in Fig. 16.
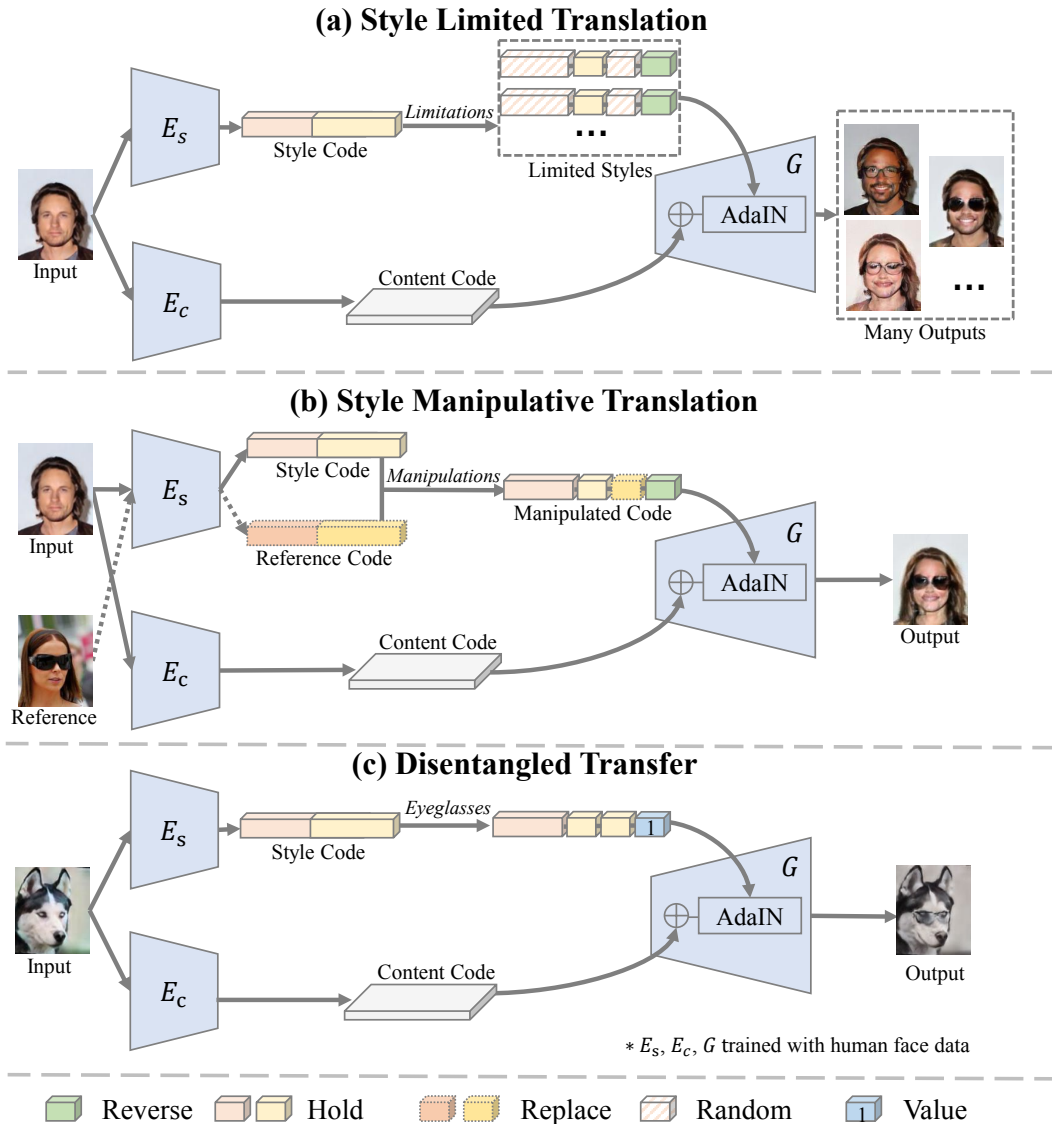


Figure 12. Detailed description of the inference phase of AGUIT. (a) Style code of input image is operated by Hold, Reversed or Random to get a variety of outputs. (b) Style code of input image is operated by Hold, Reversed or replaced by reference's style code to get the manipulated output. (c) Style code of a dog is valued to +1 at the specific dimension for controlling the presence of eyeglasses.

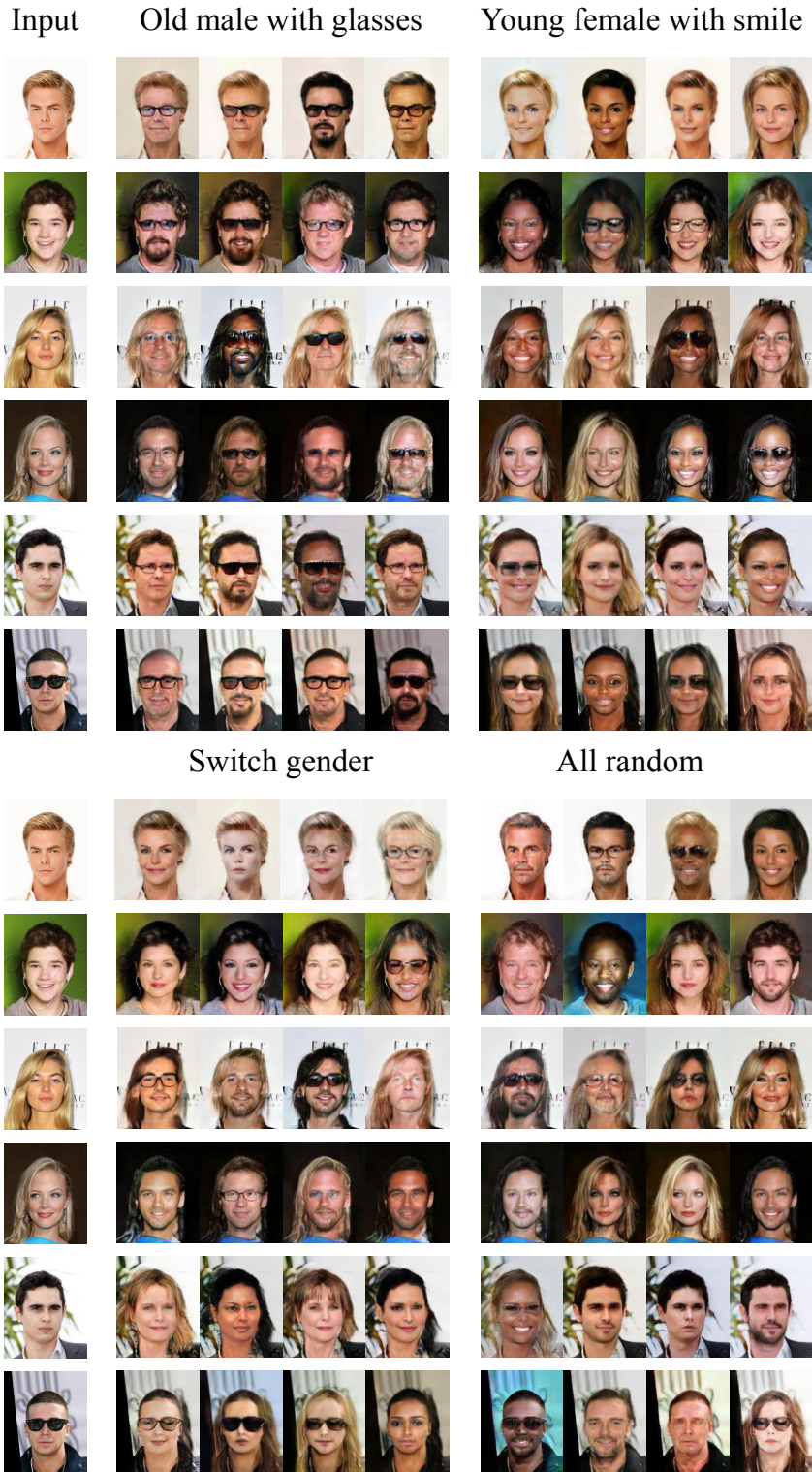| Input | Old male with glasses | Young female with smile |
|---|---|---|

| Switch gender | All random |
|---|---|

Figure 13. Additional results of AGUIT for style operable image translaiton tasks in CelebA dataset.

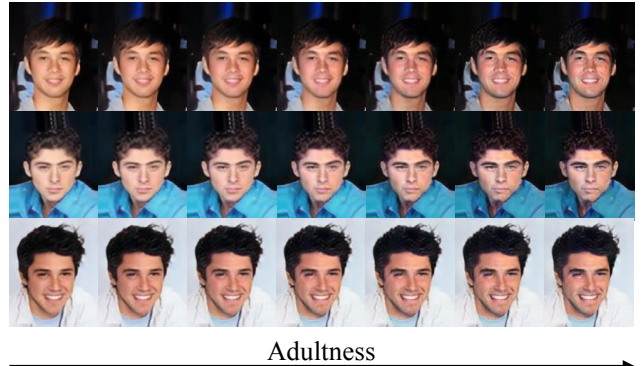Gender



Race



Young



Adultness
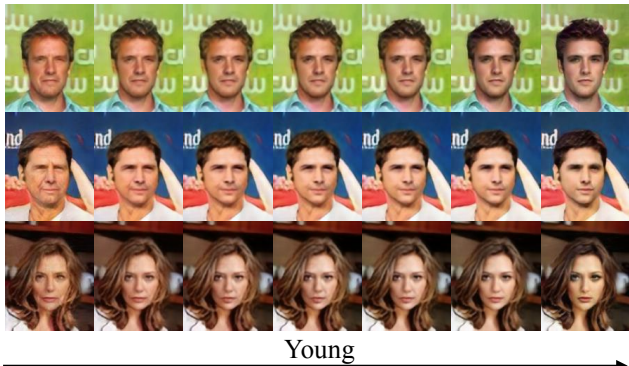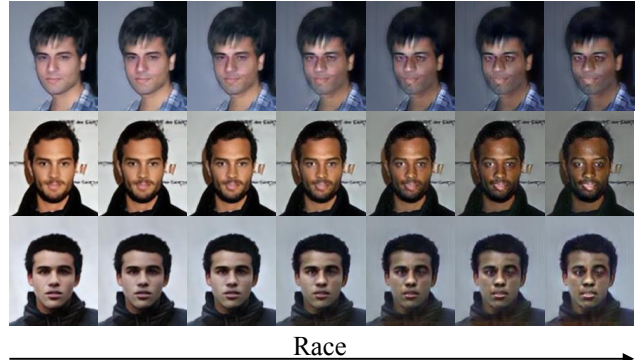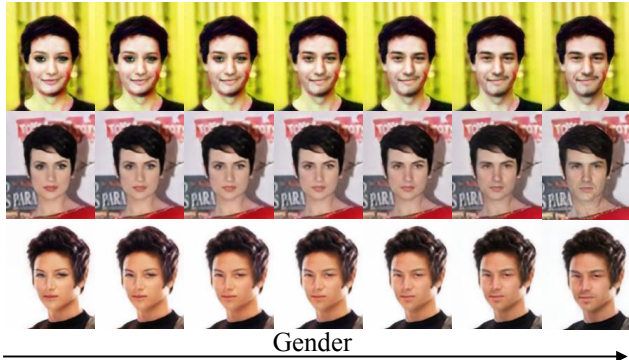


Smiling



Eyeglasses Style

Figure 14. Supervised disentanglement on the attribute part of style code.

Figure 15. Unsupervised disentanglement on the noise part of style code.

Figure 16. The results of 512x512 photos on CelebA dataset. First column is the input. More and more labels is reversed in the next two columns. The forth column is the multi-modal results. Although limited batch size causes the unsupervised disentanglement of style code not as well as 128x128 models, the training of AGUIT does not need large-batch training like ProGAN, BigGAN, and StyleGAN.