

TOWARDS IDENTIFIABLE UNSUPERVISED DOMAIN TRANSLATION: A DIVERSIFIED DISTRIBUTION MATCHING APPROACH

Sagar Shrestha & Xiao Fu *

School of Electrical Engineering and Computer Science
Oregon State University
Corvallis, OR 97331, USA
{shressag, xiao.fu}@oregonstate.edu

ABSTRACT

Unsupervised domain translation (UDT) aims to find functions that convert samples from one domain (e.g., sketches) to another domain (e.g., photos) without changing the high-level semantic meaning (also referred to as “content”). The translation functions are often sought by probability distribution matching of the transformed source domain and target domain. CycleGAN stands as arguably the most representative approach among this line of work. However, it was noticed in the literature that CycleGAN and variants could fail to identify the desired translation functions and produce content-misaligned translations. This limitation arises due to the presence of multiple translation functions—referred to as “measure-preserving automorphism” (MPA)—in the solution space of the learning criteria. Despite awareness of such identifiability issues, solutions have remained elusive. This study delves into the core identifiability inquiry and introduces an MPA elimination theory. Our analysis shows that MPA is unlikely to exist, if multiple pairs of diverse cross-domain conditional distributions are matched by the learning function. Our theory leads to a UDT learner using distribution matching over auxiliary variable-induced subsets of the domains—other than over the entire data domains as in the classical approaches. The proposed framework is the first to rigorously establish translation identifiability under reasonable UDT settings, to our best knowledge. Experiments corroborate with our theoretical claims.

1 INTRODUCTION

Domain translation (DT) aims to convert data samples from one feature domain to another, while keeping the key content information. DT naturally arises in many applications, e.g., transfer learning (Zhuang et al., 2020), domain adaptation (Ganin et al., 2016; Courty et al., 2017), and cross-domain retrieval (Huang et al., 2015). Among them, a premier application is image-to-image (I2I) translation (e.g., profile photo to cartoonized emoji and satellite images to street map plots (Isola et al., 2017)). *Supervised* domain translation (SDT) relies on paired data from the source and target domains. There, the translation functions are learned via matching the sample pairs.

Nonetheless, paired data are not always available. In *unsupervised domain translation* (UDT), the arguably most widely adopted idea is to find neural transformation functions that perform probability distribution matching of the domains. The idea emerged in the literature in early works, e.g., (Liu & Tuzel, 2016; Taigman et al., 2017; Kim et al., 2017). High-resolution image translation using distribution matching was later realized by the seminal work, namely, CycleGAN (Zhu et al., 2017). CycleGAN learns a pair of transformations that are inverse of each other. One of transformations maps the source domain to match the distribution of the target domain, and the other transformation does the opposite. The distribution matching part is realized by the generative adversarial network (GAN) (Goodfellow et al., 2014). Using GAN-based distribution matching for UDT has attracted much attention—many follow-up works emerged; see the survey (Pang et al., 2021).

*Source code is available at <https://github.com/XiaoFuLab/Identifiable-UDT.git>

Challenge - Lack of Translation Identifiability. While UDT approaches have demonstrated significant empirical success, the theoretical question of translation identifiability has received relatively limited attention. Recent works (Galanti et al., 2018b;a; Moriakov et al., 2020; Galanti et al., 2021) pointed out failure cases of CycleGAN (e.g., content-misaligned translations like those in Fig. 1) largely attribute to the lack of translation identifiability. That is, translation functions in the solution space of CycleGAN (or any distribution matching-based learners) is non-unique, due to the existence of *measure-preserving automorphism* (MPA) (Moriakov et al., 2020) (the same concept was called *density-preserving mappings* in (Galanti et al., 2018b;a)). MPA can “swap” the cross-domain sample correspondences without changing the data distribution—which is likely the main source of producing content misaligned samples after translation as seen in Fig. 1. Many efforts were made to empirically enhance the performance of UDT, via implicitly or explicitly promoting solution uniqueness of their loss functions (Liu et al., 2017; Courty et al., 2017; Xu et al., 2022; Yang et al., 2023). A number of notable works approached the identifiability/uniqueness challenge by assuming that the desired translation functions have simple (e.g., linear (Gulrajani & Hashimoto, 2022)) or specific structures (de Bézenac et al., 2021). However, translation identifiability without using such restrictive structural assumptions have remained elusive.

Contributions. In this work, we revisit distribution matching-based UDT. Our contribution lies in both identifiability theory and implementation:

- **Theory Development: Establishing Translation Identifiability.** We delve into the core theoretical challenge regarding identifiability of the translation functions. As mentioned, the solution space of existing distribution matching criteria could be easily affected by MPA. However, our analysis shows that the chance of having MPA decreases quickly when the translation function aligns more than one pair of diverse distributions. This insight allows us to come up a sufficient condition, namely, the *sufficiently diverse condition* (SDC), to establish translation identifiability of UDT. To our best knowledge, our result stands as the first UDT identifiability theory without using simplified structural assumptions.

- **Simple Implementation via Auxiliary Variables.** Our theoretical revelation naturally gives rise to a novel UDT learning criterion. This criterion aligns multiple pairs of conditional distributions across the source and target domains. We define these conditional distributions over (overlapping) sub-domains of the source/target domains using auxiliary variables. We demonstrate that in practical applications such as unpaired I2I translation, obtaining these sub-domains can be a straightforward task, e.g., through available side information or querying the foundation models like CLIP (Radford et al., 2021). Consequently, our identification theory can be readily put into practice.

Notation. The full list of notations is in the supplementary material. Notably, we use \mathbb{P}_x and $\mathbb{P}_{x|u}$ to denote the *probability measures* of x and x conditioned on u , respectively. We denote the corresponding *probability density function* (PDF) of x by $p(x)$. For a measurable function $f: \mathcal{X} \rightarrow \mathcal{Y}$ and a distribution \mathbb{P}_x defined over space \mathcal{X} , the notation $f_{\#\mathbb{P}_x}$ denotes the *push-forward measure*; that is, for any measurable set $\mathcal{A} \subseteq \mathcal{Y}$, $f_{\#\mathbb{P}_x}[\mathcal{A}] = \mathbb{P}_x[f^{\text{preimg}}(\mathcal{A})]$, where $f^{\text{preimg}}(\mathcal{A}) = \{x \in \mathcal{X} \mid f(x) \in \mathcal{A}\}$. Simply speaking, $f_{\#\mathbb{P}_x}$ denotes the distribution of $f(x)$ where $x \sim \mathbb{P}_x$. The notation $f_{\#\mathbb{P}_x} = \mathbb{P}_y$ means that the PDFs of $f(x)$ and y are identical *almost everywhere* (a.e.).

2 PRELIMINARIES

Considers two data domains (e.g., photos and sketches). The samples from the two domains are represented by $x \in \mathcal{X} \subseteq \mathbb{R}^{D_x}$ and $y \in \mathcal{Y} \subseteq \mathbb{R}^{D_y}$. We make the following assumption:

Assumption 1. For every $x \in \mathcal{X}$, it has a corresponding $y \in \mathcal{Y}$, and vice versa. In addition, there exist deterministic continuous functions $f^*: \mathcal{Y} \rightarrow \mathcal{X}$ and $g^*: \mathcal{X} \rightarrow \mathcal{Y}$ that link the corresponding pairs; i.e.,

$$f^*(y) = x, \quad g^*(x) = y, \quad \forall \text{ corresponding pair } (x, y). \quad (1)$$

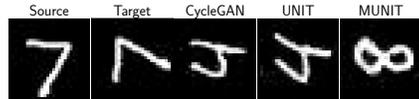


Figure 1: Lack of translation identifiability often leads to *content misalignment* in distribution matching based UDT methods, e.g., CycleGAN (Zhu et al., 2017), MUNIT (Huang et al., 2018), and UNIT (Liu et al., 2017). Source domain: MNIST Digits. Target Domain: Rotated Display of MNIST.

In the context of domain translation, a linked (x, y) pair can be regarded as cross-domain data samples that represent the same “content”, and the translation functions (f^*, g^*) are responsible for changing their “appearances/styles”. The term “content” refers to the semantic information to be kept across domains after translation. In Fig. 1, the content is the identity of the digit (other than writing style or the rotation); in Fig. 4 of Sec. 3, the content can be understood as the shared characteristics of the person in both the cartoon and the photo domains, which can collectively identify the person.

Note that in the above setting, the goal is to find *two* ground-truth translation functions where one function’s source is the other’s target. Hence, both \mathcal{X} and \mathcal{Y} can serve as the source/target domains. In addition, the above also implies $f^* = (g^*)^{-1}$, i.e., the ground-truth translation functions are invertible. Under this setting, if one can identify g^* and f^* , then the samples in one domain can be translated to the other domain—while not changing the content. Note that Assumption 1 means that there is one-to-one correspondence between samples in the two domains, which can be a somewhat stringent condition in some cases. However, as we will explain in detail later, many UDT works, e.g., CycleGAN (Zhu et al., 2017) and variants (Liu et al., 2017; Kim et al., 2017; Choi et al., 2018; Park et al., 2020), essentially used the model in Assumption 1 to attain quite interesting empirical results. This makes it a useful model and intrigues us to understand its underlying properties.

Supervised Domain Translation (SDT). In SDT, the corresponding pairs (x, y) are assumed to be aligned *a priori*. Then, learning a translation function is essentially a regression problem—e.g., via finding g (or f) such that $D(g(x)||y)$ (or $D(f(y)||x)$) is minimized over all given pairs, where $D(\cdot||\cdot)$ is a certain “distance” measure; see, e.g., (Isola et al., 2017; Wang et al., 2018).

Unsupervised Domain Translation (UDT). In UDT, samples from the two domains are acquired separately without alignment. Hence, sample-level matching as often done in SDT is not viable. Instead, UDT is often formulated as a probability distribution matching problem (see, e.g., (Zhu et al., 2017; Taigman et al., 2017; Kim et al., 2020; Park et al., 2020))—as distribution matching can be attained without using sample-level correspondences. Assume that x and y are the random vectors that represent the data from the \mathcal{X} -domain and the \mathcal{Y} -domain, respectively. Then, the desired f^* and g^* are sought via finding f and g such that

$$\mathbb{P}_y = g_{\#\mathbb{P}_x} \quad \text{and} \quad \mathbb{P}_x = f_{\#\mathbb{P}_y}. \quad (2)$$

The hope is that distribution matching can work as a surrogate of sample-level matching as in SDT. The arguably most representative work in UDT is CycleGAN (Zhu et al., 2017). The CycleGAN loss function is as follows:

$$\min_{f, g} \max_{d_x, d_y} \mathcal{L}_{\text{GAN}}(g, d_y, x, y) + \mathcal{L}_{\text{GAN}}(f, d_x, x, y) + \lambda \mathcal{L}_{\text{cyc}}(g, f), \quad (3)$$

where d_x and d_y represent two discriminators in domains \mathcal{X} and \mathcal{Y} , respectively,

$$\mathcal{L}_{\text{GAN}}(g, d_y, x, y) = \mathbb{E}_{y \sim \mathbb{P}_y} [\log d_y(y)] + \mathbb{E}_{x \sim \mathbb{P}_x} [\log(1 - d_y(g(x)))], \quad (4)$$

$\mathcal{L}_{\text{GAN}}(f, d_x, x, y)$ is defined in the same way, and the cycle-consistency term is defined as

$$\mathcal{L}_{\text{cyc}}(g, f) = \mathbb{E}_{x \sim \mathbb{P}_x} [\|f(g(x)) - x\|_1] + \mathbb{E}_{y \sim \mathbb{P}_y} [\|g(f(y)) - y\|_1]. \quad (5)$$

The minimax optimization of the \mathcal{L}_{GAN} terms enforces $g_{\#\mathbb{P}_x} = \mathbb{P}_y$ and $f_{\#\mathbb{P}_y} = \mathbb{P}_x$. The \mathcal{L}_{cyc} term encourages $f = g^{-1}$. CycleGAN showed the power of distribution matching in UDT and has triggered a lot of interests in I2I translation. Many variants of CycleGAN were also proposed to improve the performance; see the survey (Pang et al., 2021).

Lack of Translation Identifiability, MPA and Content Misalignment. Many works have noticed that distribution matching-type learning criterion may suffer from the lack of translation identifiability (Liu et al., 2017; Moriakov et al., 2020; Galanti et al., 2018b; 2021; Xu et al., 2022); i.e., the solution space of these criteria could have multiple solutions, and thus lack the ability to recover the ground-truth g^* and f^* . The lack of identifiability often leads to issues such as content misalignment as we saw in Fig. 1. To understand the identifiability challenge, let us formally define identifiability of any bi-directional UDT learning criterion:

Definition 1. (*Identifiability*) Under the setting of Assumption 1, assume that (\hat{f}, \hat{g}) is any optimal solution of a UDT learning criterion. Then, identifiability of (f^*, g^*) holds under the UDT learning criterion if and only if $\hat{f} = f^*$ and $\hat{g} = g^*$ a.e.

Notice that we used the *optimal solution* in the definition. This is because identifiability is a characterization of the “kernel space” (which contains all the zero-loss solutions) of a learning criterion (Moriakov et al., 2020; Fu et al., 2019). In other words, when a UDT criterion admits translation identifiability, it indicates that the criterion provides a valid objective for the learning task—but identifiability is not related to the optimization procedure. We will also use the following:

Definition 2. (MPA) A *measure-preserving automorphism (MPA)* of \mathbb{P}_x is a continuous function $h : \mathcal{X} \rightarrow \mathcal{X}$ such that $\mathbb{P}_x = h_{\#}\mathbb{P}_x$.

Simply speaking, MPA defined in this work is the continuous transformation $h(x)$ whose output has the same PDF as $p(x)$. Take the one-dimensional Gaussian distribution $x \sim \mathcal{N}(\mu, \sigma^2)$ as an example. The MPA of $\mathcal{N}(\mu, \sigma^2)$ is $h(x) = -x + 2\mu$. A recent work (Moriakov et al., 2020) suggested that non-identifiability of the desired translation functions by CycleGAN is caused by the existence of MPA. Their finding can be summarized in the following Fact:

Fact 1. If MPA of \mathbb{P}_x or \mathbb{P}_y exists, then CycleGAN and any criterion using distribution matching in (2) do not have identifiability of f^* and g^* .

Proof: It is straightforward to see that $\mathbb{P}_y = g_{\#}\mathbb{P}_x$ and $\mathbb{P}_x = f_{\#}\mathbb{P}_y$. In addition, f^* and g^* are invertible. Hence, the ground truth (f^*, g^*) is an optimal solution of CycleGAN that makes the loss in (3) equal to zero. However, due to the existence MPA, one can see that $\hat{f} = h \circ f^*$ can also attain $\mathbb{P}_x = \hat{f}_{\#}\mathbb{P}_y$. This is because we have $\hat{f}_{\#}\mathbb{P}_y = h \circ f_{\#}\mathbb{P}_y = h_{\#}\mathbb{P}_x = \mathbb{P}_x$.

Plus, as $h \circ f^*$ is still invertible, \hat{f} still makes the cycle-consistency loss zero. Hence, the solution of CycleGAN is not unique and this loses identifiability of the ground truth translation functions. \square

The existence of MPA in the solution space of the UDT learning losses may be detrimental in terms of avoiding content misalignment. To see this, consider the example in Fig. 2. There, $\mathbb{P}_x = \mathcal{N}(\mu, \sigma^2)$ and $h(x) = -x + 2\mu$ is an MPA of \mathbb{P}_x , as mentioned. Note that $\hat{f} = h \circ f^*$ can be an optimal solution found by CycleGAN. However, such an \hat{f} can cause misalignment. To explain, assume $x = a$ and $y = b$ are associated with the same entity, which means that $a = f^*(b)$ represents the ground-truth alignment and translation. However, as $p(-a + 2\mu) = p(h(a)) = p(h \circ f^*(b)) = p(\hat{f}(b))$, the learned function \hat{f} wrongly translates $y = b$ to $x = -a + 2\mu$.

Our Gaussian example seems to be special as it has symmetry about its mean. However, the existence of MPA is not unusual. To see this, we show the following result:

Proposition 1. Suppose that \mathbb{P}_x admits a continuous PDF, $p(x)$ and $p(x) > 0, \forall x \in \mathcal{X}$. Assume that \mathcal{X} is simply connected. Then, there exists a continuous non-trivial (non-identity) $h(\cdot)$ such that $h_{\#}\mathbb{P}_x = \mathbb{P}_x$.

Note that there are similar results in (Moriakov et al., 2020) regarding the existence of MPA, but more assumptions were made in their proof. The universal existence of MPA attests to the challenging nature of establishing translation identifiability in UDT.

3 IDENTIFIABLE UDT VIA DIVERSIFIED DISTRIBUTION MATCHING

Intuition - Exploiting Diversity of Distributions. Our idea starts with the following observation: If two distributions have different PDFs, a shared MPA is unlikely to exist. Fig. 3 illustrates the

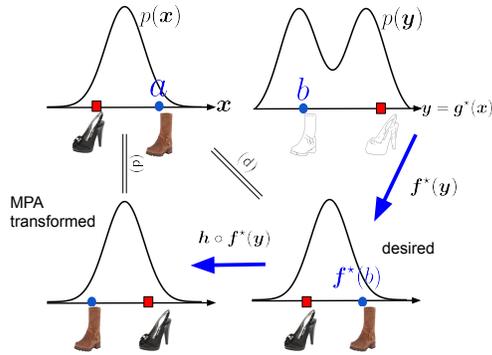


Figure 2: Illustration of the lack of identifiability and MPA-induced content misalignment; “ $\stackrel{(d)}{=}$ ” means distribution matching.

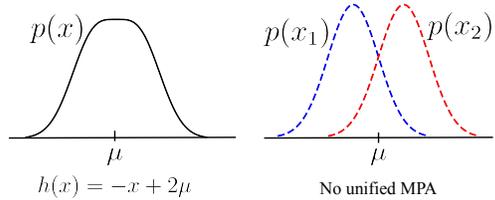


Figure 3: A unified MPA is harder to exist for a group of distributions.

intuition. Consider two Gaussian distributions $x_1 \sim \mathcal{N}(\mu_1, 1)$ and $x_2 \sim \mathcal{N}(\mu_2, 1)$ with $\mu_1 \neq \mu_2$. For each of them, $h(x) = -x + 2\mu_i$ for $i = 1, 2$ is an MPA. However, there is not a function that can serve as a unified MPA to attain $h_{\#\mathbb{P}_{x_1}} = \mathbb{P}_{x_1}$ & $h_{\#\mathbb{P}_{x_2}} = \mathbb{P}_{x_2}$ simultaneously. Intuitively, the diversity of the PDFs of x_1 and x_2 has made finding a unified MPA $h(\cdot)$ difficult. This suggests that instead of matching the distributions of \mathbf{x} and $\mathbf{f}(\mathbf{y})$ and those of \mathbf{y} and $\mathbf{g}(\mathbf{x})$, it may be beneficial to match the distributions of more variable pairs whose probability measures are diverse.

Auxiliary Variable-Assisted Distribution Diversification. In applications, the corresponding samples \mathbf{x}, \mathbf{y} often share some aspects/traits. For example, in Fig. 4, the corresponding \mathbf{x} and \mathbf{y} both have dark hair or the same gender. If we model a collection of such traits as different realizations of discrete random variable u , the alphabet of u , denoted as $\{u_1, \dots, u_I\}$ represents these traits. We should emphasize that the traits is a result of the desired content invariance across domains, but need not to represent the whole content.

To proceed, we observe that the conditional distributions $\mathbb{P}_{\mathbf{x}|u=u_i}$ and $\mathbb{P}_{\mathbf{y}|u=u_i}$ satisfy $\mathbb{P}_{\mathbf{x}|u=u_i} = \mathbf{f}_{\#\mathbb{P}_{\mathbf{y}|u=u_i}}^*$, $\mathbb{P}_{\mathbf{y}|u=u_i} = \mathbf{g}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}}^*$, $\forall i$. The above holds since \mathbf{x} and \mathbf{y} have a deterministic relation and because the trait u_i is shared by the content-aligned pairs (\mathbf{x}, \mathbf{y}) .

In practice, u can take various forms. In I2I translation, one may use image categories or labels, if available, to serve as u . Note that knowing the image categories does *not* mean the samples from the two domains are aligned, as each category could contain a large amount of samples. In addition, one can use sample attributes (such as hair color, gender as in Fig. 4) to serve as u , if these attributes are not meant to be changed in the considered translation tasks. If not immediately available, these attributes can be annotated by open-sourced AI models, e.g., CLIP (Radford et al., 2021); see detailed implementation in the supplementary material. A similar idea of using CLIP to acquire auxiliary information was explored in (Gabbay et al., 2021).

By Proposition 1, it is almost certain that $\mathbb{P}_{\mathbf{x}|u=u_i}$ has an MPA h_i for all $i \in [I]$. However, it is likely that $h_i \neq h_j$ if $\mathbb{P}_{\mathbf{x}|u=u_i}$ and $\mathbb{P}_{\mathbf{x}|u=u_j}$ are sufficiently different. As a consequence, similar to what we saw in Fig. 3, if one looks for \mathbf{f} that does simultaneous matching of

$$\mathbb{P}_{\mathbf{x}|u=u_i} = \mathbf{f}_{\#\mathbb{P}_{\mathbf{y}|u=u_i}}, \forall i \in [I], \quad (6)$$

it is more possible that $\mathbf{f} = \mathbf{f}^*$ instead of having other solutions—this leads to identifiability of \mathbf{f}^* .

Proposed Loss Function. We propose to match multiple distribution pairs $(\mathbb{P}_{\mathbf{x}|u_i}, \mathbf{f}_{\#\mathbb{P}_{\mathbf{y}|u_i}})$ (as well as $(\mathbb{P}_{\mathbf{y}|u_i}, \mathbf{g}_{\#\mathbb{P}_{\mathbf{x}|u_i}})$) for $i = 1, \dots, I$. For each pair, we use discriminator $\mathbf{d}_x^{(i)} : \mathcal{X} \rightarrow [0, 1]$ (and $\mathbf{d}_y^{(i)} : \mathcal{Y} \rightarrow [0, 1]$ in reverse direction). Then, our loss function is as follows:

$$\min_{\mathbf{f}, \mathbf{g}} \max_{\{\mathbf{d}_x^{(i)}, \mathbf{d}_y^{(i)}\}} \sum_{i=1}^I \left(\mathcal{L}_{\text{GAN}}(\mathbf{g}, \mathbf{d}_y^{(i)}, \mathbf{x}, \mathbf{y}) + \mathcal{L}_{\text{GAN}}(\mathbf{f}, \mathbf{d}_x^{(i)}, \mathbf{x}, \mathbf{y}) \right) + \lambda \mathcal{L}_{\text{cyc}}(\mathbf{g}, \mathbf{f}), \quad (7)$$

where we have

$$\mathcal{L}_{\text{GAN}}(\mathbf{g}, \mathbf{d}_y^{(i)}, \mathbf{x}, \mathbf{y}) = \Pr(u = u_i) \left(\mathbb{E}_{\mathbf{y} \sim \mathbb{P}_{\mathbf{y}|u=u_i}} \left[\log \mathbf{d}_y^{(i)}(\mathbf{y}) \right] + \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\mathbf{x}|u=u_i}} \left[\log \left(1 - \mathbf{d}_y^{(i)}(\mathbf{g}(\mathbf{x})) \right) \right] \right).$$

Note that $\mathbf{x} \sim \mathbb{P}_{\mathbf{x}|u_i}$ represents samples that share the same characteristic defined by u_i (e.g., hair color, eye color, gender). This means that the loss function matches a suite of distributions defined over (potentially overlapping) subdomains over the entire domain \mathcal{X} and \mathcal{Y} . We should emphasize that the auxiliary variable is only needed in the training stage, but not the testing stage.

We call the proposed method *diversified distribution matching for unsupervised domain translation* (DIMENSION)¹. The following lemma shows that DIMENSION exactly realizes our idea in (6):

Lemma 1. Assume that an optimal solution of (7) is $(\hat{\mathbf{f}}, \hat{\mathbf{g}}, \{\hat{\mathbf{d}}_x^{(i)}, \hat{\mathbf{d}}_y^{(i)}\})$. Then, under Assumption 1, we have $\mathbb{P}_{\mathbf{x}|u=u_i} = \hat{\mathbf{f}}_{\#\mathbb{P}_{\mathbf{y}|u=u_i}}$, $\mathbb{P}_{\mathbf{y}|u=u_i} = \hat{\mathbf{g}}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}}$, $\forall i \in [I]$, and $\hat{\mathbf{f}} = \hat{\mathbf{g}}^{-1}$, a.e.

¹Note that we still use the term “unsupervised” despite the need of auxiliary information—as no paired samples are required. We avoided using “semi-supervised” or “weakly supervised” as these are often reserved for methods using some paired samples; see, e.g., (Wang et al., 2020; Mustafa & Mantiuk, 2020).

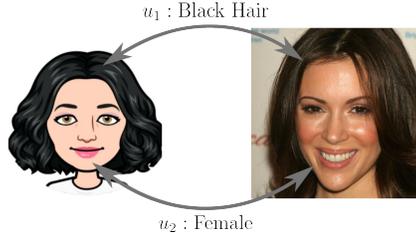


Figure 4: Examples of u_i .

Identifiability Characterization. Lemma 1 means that solving the DIMENSION loss leads to conditional distribution matching as we hoped for in (6). However, it does not guarantee that $(\hat{\mathbf{f}}, \hat{\mathbf{g}})$ found by DIMENSION satisfies $\hat{\mathbf{f}} = \mathbf{f}^*$ and $\hat{\mathbf{g}} = \mathbf{g}^*$. Towards establishing *identifiability* of the ground-truth translation functions via DIMENSION, we will use the following definition:

Definition 3 (Admissible MPA). *Given auxiliary variable u , the function $\mathbf{h}(\cdot)$ is said to be an admissible MPA of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ if and only if $\mathbb{P}_{\mathbf{x}|u=u_i} = \mathbf{h}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}}, \forall i \in [I]$.*

Now, due to the deterministic relationship between the pair \mathbf{x} and \mathbf{y} , we have the following fact:

Fact 2. *Suppose that Assumption 1 holds. Then, there exists an admissible MPA of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ if and only if there exists an admissible MPA of $\{\mathbb{P}_{\mathbf{y}|u=u_i}\}_{i=1}^I$.*

The above means that if we establish that there is no admissible MPA of the $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$, it suffices to conclude that there is no admissible MPA of $\{\mathbb{P}_{\mathbf{y}|u=u_i}\}_{i=1}^I$.

As described before, to ensure identifiability of the translation functions via solving the DIMENSION loss, we hope the conditional distributions $\mathbb{P}_{\mathbf{x}|u=u_i}$ and $\mathbb{P}_{\mathbf{y}|u=u_i}$ to be sufficiently different. We formalize this requirement in the following definition:

Definition 4 (Sufficiently Diverse Condition (SDC)). *For any two disjoint sets $\mathcal{A}, \mathcal{B} \subset \mathcal{X}$, where \mathcal{A} and \mathcal{B} are connected, open, and non-empty, there exists a $u_{(\mathcal{A}, \mathcal{B})} \in \{u_1, \dots, u_I\}$ such that $\mathbb{P}_{\mathbf{x}|u=u_{(\mathcal{A}, \mathcal{B})}}[\mathcal{A}] \neq \mathbb{P}_{\mathbf{x}|u=u_{(\mathcal{A}, \mathcal{B})}}[\mathcal{B}]$. Then, the set of conditional distributions $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ is called sufficiently diverse.*

Definition 4 puts the desired ‘‘diversity’’ into context. It is important to note that the SDC only requires the *existence* of a certain $u_{(\mathcal{A}, \mathcal{B})} \in \{u_1, \dots, u_I\}$ for a given disjoint set pair $(\mathcal{A}, \mathcal{B})$. It does not require a unified u for all pairs; i.e., $u_{(\mathcal{A}, \mathcal{B})}$ needs not to be the same as $u_{(\mathcal{A}', \mathcal{B}')}$ for $(\mathcal{A}, \mathcal{B}) \neq (\mathcal{A}', \mathcal{B}')$. Fig. 5 shows a simple example where the two conditional distributions satisfy the SDC. In more general cases, this implies that if the PDFs of the conditional distributions exhibit different ‘‘shapes’’ over their supports, SDC is likely to hold. Using SDC, we show the following translation identifiability result:

Theorem 1 (Identifiability). *Suppose that Assumption 1 holds. Let $\mathbb{E}_{i,j}$ denote the event that the pair $(\mathbb{P}_{\mathbf{x}|u=u_i}, \mathbb{P}_{\mathbf{x}|u=u_j})$ does not satisfy the SDC. Assume that $\Pr[\mathbb{E}_{i,j}] \leq \rho$ for any $i \neq j$, where $i, j \in [I]$. Let $(\hat{\mathbf{f}}, \hat{\mathbf{g}})$ be from an optimal solution of the DIMENSION loss (7). Then, there is no admissible MPA of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ of the solution, i.e., $\hat{\mathbf{f}} = \mathbf{f}^*$, a.e. and $\hat{\mathbf{g}} = \mathbf{g}^*$, a.e. with a probability of at least $1 - \rho^{\binom{I}{2}}$.*

Theorem 1 shows that if the conditional distributions are sufficiently diverse, solving (7) can correctly identify the ground-truth translation functions. Theorem 1 also spells out the importance of having more u_i ’s (which means more auxiliary information). The increase of I improves the probability of success quickly.

Towards More Robust Identifiability. Theorem 1 uses the fact that the SDC holds with high probability for every pair of $(\mathbb{P}_{\mathbf{x}|u_i}, \mathbb{P}_{\mathbf{x}|u_j})$ (cf. $\Pr[\mathbb{E}_{i,j}] \leq \rho$). It is also of interest to see if the method is robust to violation of the SDC. To this end, consider the following condition:

Definition 5 (Relaxed Condition: r -SDC). *Let $\text{dia}(\mathcal{A}) = \sup_{\mathbf{w}, \mathbf{z} \in \mathcal{A}} \|\mathbf{w} - \mathbf{z}\|_2$ and $\mathcal{V}_{i,j} = \{(\mathcal{A}, \mathcal{B}) \mid \mathbb{P}_{\mathbf{x}|u_i}[\mathcal{A}] = \mathbb{P}_{\mathbf{x}|u_i}[\mathcal{B}] \ \& \ \mathbb{P}_{\mathbf{x}|u_j}[\mathcal{A}] = \mathbb{P}_{\mathbf{x}|u_j}[\mathcal{B}], \mathcal{A} \cap \mathcal{B} = \emptyset\}$, where \mathcal{A}, \mathcal{B} are non-empty, open and connected. Denote $M_{i,j} = \max_{(\mathcal{A}, \mathcal{B}) \in \mathcal{V}_{i,j}} \max\{\text{dia}(\mathcal{A}), \text{dia}(\mathcal{B})\}$. Then, $(\mathbb{P}_{\mathbf{x}|u_i}, \mathbb{P}_{\mathbf{x}|u_j})$ satisfies the r -SDC if $M_{i,j} \leq r$ for $r \geq 0$.*

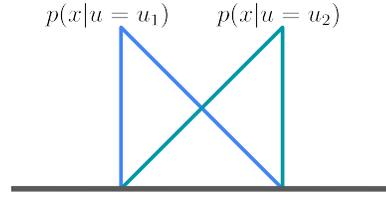


Figure 5: Conditional PDFs $p(\mathbf{x}|u = u_1)$ and $p(\mathbf{x}|u = u_2)$ that satisfy the SDC.

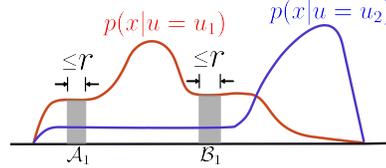


Figure 6: Illustration of relaxed SDC (r -SDC).

Note that the r -SDC becomes the SDC when $r = 0$. Unlike SDC in Definition 4, the relaxed SDC condition allows the violation of SDC over regions $\mathcal{V}_{i,j}$. Our next theorem shows that the translation identifiability still approximately holds, as long as the largest region in $\mathcal{V}_{i,j}$ is not substantial:

Theorem 2 (Robust Identifiability). *Suppose that Assumption 1 holds with \mathbf{g}^* being L -Lipschitz continuous, and that any pair of $(\mathbb{P}_{\mathbf{x}|u_i}, \mathbb{P}_{\mathbf{x}|u_j})$ satisfies the r -SDC (cf. Definition 5) with probability at least $1 - \gamma$, i.e., $\Pr[M_{i,j} \geq r] \leq \gamma$ for any $i \neq j$, where $(i, j) \in [I] \times [J]$. Let $\hat{\mathbf{g}}$ be from any optimal solution of the DIMENSION loss in (7). Then, we have $\|\hat{\mathbf{g}}(\mathbf{x}) - \mathbf{g}^*(\mathbf{x})\|_2 \leq 2rL$, $\forall \mathbf{x} \in \mathcal{X}$, with a probability of at least $1 - \gamma^{\binom{I}{2}}$. The same holds for $\hat{\mathbf{f}}$.*

Theorem 2 asserts that the estimation error of $\hat{\mathbf{g}}$ scales linearly with the “degree” of violation of the SDC (measured by r). The result is encouraging: It shows that even if the SDC is violated, the performance of DIMENSION will not decline drastically. The Lipschitz continuity assumption in Theorem 2 is mild. Note that translation functions are often represented by neural networks in practice, and neural networks with bounded weights are Lipschitz continuous functions (Bartlett et al., 2017). Hence, the numerical successes of many neural UDT models (e.g., CycleGAN) suggest that assuming that Lipschitz continuous ground-truth translation functions exist is reasonable.

4 RELATED WORKS

Prior to CycleGAN (Zhu et al., 2017), the early works (Liu & Tuzel, 2016; Taigman et al., 2017; Kim et al., 2017) started using GAN-based neural structures for distribution matching in the context of I2I translation. Similar ideas appeared in UDT problems in NLP (e.g., machine translation) (Conneau et al., 2017; Lample et al., 2017). In the literature, it was noticed that distribution matching modules lack solution uniqueness, and many works proposed remedies (see, e.g., (Liu et al., 2017; Xu et al., 2022; Xie et al., 2022; Park et al., 2020)). These approaches have worked to various extents empirically, but the translation identifiability question was unanswered. The term “content” was used in the vision literature (in the context of I2I translation) to refer to domain-invariant attributes (e.g., pose and orientation (Kim et al., 2020; Amodio & Krishnaswamy, 2019; Wu et al., 2019; Yang et al., 2023)). This is a narrower interpretation of content relative to ours—as content in our case can be high-level or latent semantic meaning that is not represented by specific attributes. Our definition of content is closer to that in multimodal and self-supervised learning (Von Kügelgen et al., 2021; Lyu et al., 2022; Daunhawer et al., 2023). Before our work, auxiliary information was also considered in UDT. For example, semi-supervised UDT (see, e.g., (Wang et al., 2020; Mustafa & Mantiuk, 2020)) uses a small set of paired data samples, but our method does not use any sample-level pairing information. Attribute-guided I2I translation (see, e.g., (Li et al., 2019; Choi et al., 2018; 2020)) specifies the desired attributes in the target domain to “guide” the translation. These are different from our auxiliary variables that can be both sample attributes or high-level concepts (which is closer to the “auxiliary variables” in nonlinear independent component analysis works, e.g., (Hyvarinen et al., 2019)). Again, translation identifiability was not considered for semi-supervised or attribute-guided UDT. There has been efforts towards understanding the translation identifiability of CycleGAN. The works of Galanti et al. (2018b;a) recognized that the success of UDT may attribute to the existence of a small number of MPAs. Moriakov et al. (2020) showed that MPA exists in the solution space of CycleGAN, and used it to explain the ill-posedness of CycleGAN. Chakrabarty & Das (2022) studied the finite sample complexity of CycleGAN in terms of distribution matching and cycle consistency. Gulrajani & Hashimoto (2022) and de Bézenac et al. (2021) argued that if the target translation functions have known structures (e.g., linear or optimal transport structures), then translation identifiability can be established. However, these conditions can be restrictive. Translation identifiability without using such structural assumptions had remained unclear before our work.

5 NUMERICAL VALIDATION

Constructing Challenging Translation Tasks. We construct challenging translation tasks to validate our theorems and to illustrate the importance of translation identifiability. To this end, we make three datasets. The first two are “MNIST v.s. Rotated MNIST” (MrM) and “Edges v.s. Rotated Shoes” (ErS). In both datasets, the rotated domains consist of samples from the “MNIST” and “Shoes” with a 90 degree rotation, respectively. We intentionally make this rotation, as rotation is

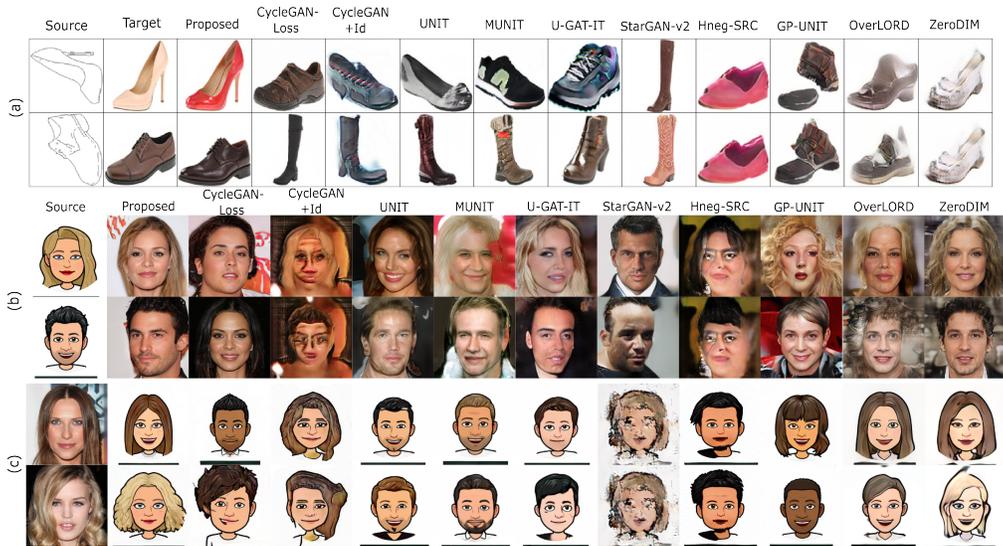


Figure 8: Qualitative results on (a) Edges to Rotated Shoes, (b) Bitmoji Faces to CelebA-HQ, and (c) CelebA-HQ to Bitmoji Faces tasks. More comprehensive illustrations are in the appendix.

a large geometric change across domains. This type of large geometric change poses a challenging translation task (Kim et al., 2020; Wu et al., 2019; Amodio & Krishnaswamy, 2019; Yang et al., 2023). In addition, we construct a task ‘‘CelebA-HQ (Karras et al., 2017) v.s. Bitmoji (Mozafari, 2020)’’ (CB). In this task, profile photos of celebrities are translated to cartoonized bitmoji figures, and vice versa. We intentionally choose these two domains to make the translation challenging: The profile photos have rich details and are diverse in terms of face orientation, expression, hair style, etc., but the Bitmoji pictures have a relatively small set of choices of these attributes (e.g., they are always front-facing). More details of the datasets are in Sec. F.4 in the supplementary material.

Baselines. The baselines include some representative UDT methods and some recent developments, i.e., GP-UNIT (Yang et al., 2023), Hneg-SRC (Jung et al., 2022), OverLORD (Gabbay & Hoshen, 2021), ZeroDIM (Gabbay et al., 2021), StarGAN-v2 (Choi et al., 2020), U-GAT-IT (Kim et al., 2020), MUNIT (Huang et al., 2018), UNIT (Liu et al., 2017), and CycleGAN (Zhu et al., 2017). In particular, two versions of CycleGAN are used. ‘‘CycleGAN Loss’’ refers to the plan-vanilla CycleGAN objective in (3) and CycleGAN+Id refers to the ‘‘identity-regularized’’ version in (Zhu et al., 2017). ZeroDIM uses the same auxiliary information as that used by the proposed method.

MNIST to Rotated MNIST. Fig. 7 shows the results. In this case, we use $u \in \{1, \dots, 10\}$, i.e., the labels of the identity of digits, as the alphabet of the auxiliary variable. Note that knowing such labels does not mean that the cross-domain pairs (x, y) are known. Alternatively, one can also use digit shapes as the alphabets (see Sec. F.6). One can see that DIMENSION learns to translate the digits to their corresponding rotated versions. But the baselines sometimes misalign the samples. The results are consistent with our analysis (see Sec. F.6 for more results).

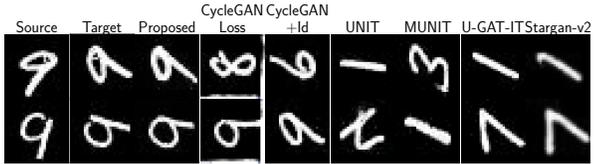


Figure 7: Translation from MNIST to rotated MNIST.

Edges to Rotated Shoes. From Fig. 8 (a), one can see that the baselines all misalign the edges with wrong shoes. Instead, the proposed DIMENSION, using the shoe types (shoes, boots, sandals, and slippers) as the alphabet of u , does not encounter this issue. More experiments including the reverse translation (i.e., shoes to edges) are in Sec. F.6 in the supplementary material.

CelebA-HQ and Bitmoji. Figs. 8 (b-c) show the results. The proposed method uses $u \in \{‘male’, ‘female’, ‘black hair’, ‘non-black hair’\}$. To obtain the auxiliary information for each sample, we use CLIP to automatically annotate the images. A remark is that translating from the Bitmoji domain to the CelebA-HQ domain [see. Fig. 8 (b)] is particularly hard. This is because the learned translation function needs to ‘‘fill in’’ a lot of details to make the generated profiles

Table 1: LPIPS scores for the ErS and MrM tasks and FID scores for all tasks. E: Edges, rS: rotated Shoes, M: MNIST, rM: rotated MNIST, C: CelebA-HQ, B: Bitmoji faces.

Method	LPIPS (\downarrow)				FID (\downarrow)					
	E \rightarrow rS	rS \rightarrow E	M \rightarrow rM	rM \rightarrow M	E	rS	M	rM	C	B
Proposed	0.29 \pm 0.06	0.35 \pm 0.10	0.11 \pm 0.08	0.09 \pm 0.04	21.47	40.14	13.95	16.07	32.03	20.50
CycleGAN-Loss	0.43 \pm 0.06	0.50 \pm 0.07	0.34 \pm 0.07	0.33 \pm 0.09	35.83	55.42	16.09	16.11	36.71	28.02
CycleGAN	0.65 \pm 0.03	0.54 \pm 0.07	0.27 \pm 0.09	0.28 \pm 0.09	259.31	130.84	46.05	34.01	196.52	85.05
U-GAT-IT	0.56 \pm 0.05	0.48 \pm 0.07	0.25 \pm 0.09	0.25 \pm 0.09	288.03	58.20	11.78	11.67	50.28	39.09
UNIT	0.49 \pm 0.03	0.58 \pm 0.03	0.25 \pm 0.06	0.25 \pm 0.08	33.95	96.28	20.44	19.15	53.63	33.56
MUNIT	0.50 \pm 0.03	0.58 \pm 0.04	0.28 \pm 0.09	0.28 \pm 0.09	43.83	86.68	14.89	15.96	62.49	27.59
StarGAN-v2	0.39 \pm 0.05	0.52 \pm 0.11	0.28 \pm 0.09	0.29 \pm 0.10	75.46	138.34	30.07	32.20	35.44	282.98
Hneg-SRC	0.45 \pm 0.06	0.50 \pm 0.07	-	-	210.27	198.77	-	-	129.34	66.36
GP-UNIT	0.49 \pm 0.08	0.44 \pm 0.05	-	-	231.31	96.32	-	-	32.40	30.30
OverLORD	0.43 \pm 0.06	0.42 \pm 0.05	-	-	101.14	124.02	-	-	76.10	31.08
ZeroDIM	0.38 \pm 0.06	0.41 \pm 0.07	-	-	85.56	187.45	-	-	88.36	36.21

“-” means that method is not applicable to the dataset due to small resolution.

photorealistic. Our method clearly outperforms the baselines in both directions of translation; see more in Sec. F.6 in the supplementary material.

Metrics and Quantative Evaluation. We employ two widely adopted metrics in UDT. The first is the *learned perceptual image patch similarity* (LPIPS) (Zhang et al., 2018), which leverages the known ground-truth correspondence between (x, y) . LPIPS measures the “perceptual distance” between the translated images and the ground-truth target images. In addition, we also use the *Fréchet inception distance* (FID) score (Heusel et al., 2017) in all tasks. FID measures the visual quality of the learned translation using a distribution divergence between the translated images and the target domain. In short, LPIPS and FID correspond to the content alignment performance and the target domain-attaining ability, respectively; see details of the metrics Sec. F.4.

Table 1 shows the LPIPS scores over the first two datasets where the ground-truth pairs are known. One can see that DIMENSION significantly outperforms the baselines—which is a result of good content alignment. The FID scores in the same table show that our method produces translated images that have similar characteristics of the target domains. The FID scores output by our method are either the lowest or the second lowest.

Detailed Settings and More Experiments. See Sec. E-H for settings and more results.

6 CONCLUSION

In this work, we revisited the UDT and took a deep look at a core theoretical challenge, namely, the translation identifiability issue. Existing UDT approaches (such as CycleGAN) often lack translation identifiability and may produce content-misaligned translations. This issue largely attributes to the presence of MPA in the solution space of their distribution matching modules. Our approach leverages the existence of domain-invariant auxiliary variables to establish translation identifiability, using a novel diversified distribution matching criterion. To our best knowledge, the identifiability result stands as the first of its kind, without using restrictive conditions on the structure of the desired translation functions. We also analyzed the robustness of proposed method when the key sufficient condition for identifiability is violated. Our identifiability theory leads to an easy-to-implement UDT system. Synthetic and real-data experiments corroborated with our theoretical findings.

Limitations. Our work considers a model where the ground-truth translation functions are deterministic and bijective. This setting has been (implicitly or explicitly) adopted by a large number of existing works, with the most notable representative being CycleGAN. However, there can be multiple “correct” translation functions in UDT, as the same “content” can be combined with various “styles”. Such cases may be modeled using probabilistic translation mechanisms (Huang et al., 2018; Choi et al., 2020; Yang et al., 2023), yet the current analytical framework needs a significant revision to accommodate the probabilistic setting. In addition, our method makes use of auxiliary variables that may be nontrivial to acquire in certain cases. We have shown that open-sourced foundation models such as CLIP can help acquire such auxiliary variables and that the method is robust to noisy/wrong auxiliary variables (see Sec. H). However, it is still of great interest to develop provable UDT translation schemes without using auxiliary variables.

Acknowledgement. This work is supported in part by the Army Research Office (ARO) under Project ARO W911NF-21-1-0227, and in part by the National Science Foundation (NSF) CAREER Award ECCS-2144889.

REFERENCES

- Matthew Amodio and Smita Krishnaswamy. TravelGAN: Image-to-image translation by transformation vector learning. In *Proceedings of IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*, pp. 8983–8992, 2019.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017.
- Neal L Carothers. *Real analysis*. Cambridge University Press, 2000.
- Anish Chakrabarty and Swagatam Das. On translation and reconstruction guarantees of the cycle-consistent generative adversarial networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:23607–23620, 2022.
- Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. StarGAN: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*, pp. 8789–8797, 2018.
- Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. StarGAN v2: Diverse image synthesis for multiple domains. In *Proceedings of IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*, pp. 8188–8197, 2020.
- Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*, 2017.
- Nicolas Courty, Rémi Flamary, Amaury Habrard, and Alain Rakotomamonjy. Joint distribution optimal transportation for domain adaptation. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017.
- George Darmais. Analyse des liaisons de probabilité. In *Proceedings of International Statistic Conferences*, pp. 231, 1951.
- Imant Daunhawer, Alice Bizeul, Emanuele Palumbo, Alexander Marx, and Julia E Vogt. Identifiability results for multimodal contrastive learning. *arXiv preprint arXiv:2303.09166*, 2023.
- Emmanuel de Bézenac, Ibrahim Ayed, and Patrick Gallinari. CycleGAN through the lens of (dynamical) optimal transport. In *Proceedings of Joint European Conference on Machine Learning and Knowledge Discovery in Databases (ECML)*, pp. 132–147. Springer, 2021.
- Xiao Fu, Kejun Huang, Nicholas D Sidiropoulos, and Wing-Kin Ma. Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications. *IEEE Signal Processing Magazine*, 36(2):59–80, 2019.
- Aviv Gabbay and Yedid Hoshen. Scaling-up disentanglement for image translation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (CVPR)*, pp. 6783–6792, 2021.
- Aviv Gabbay, Niv Cohen, and Yedid Hoshen. An image is worth more than a thousand words: Towards disentanglement in the wild. *Advances in Neural Information Processing Systems (NeurIPS)*, 34:9216–9228, 2021.
- Tomer Galanti, Sagie Benaim, and Lior Wolf. Generalization bounds for unsupervised cross-domain mapping with WGANs. *arXiv preprint arXiv:1807.08501*, 2018a.

- Tomer Galanti, Lior Wolf, and Sagie Benaim. The role of minimal complexity functions in unsupervised learning of semantic mappings. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2018b.
- Tomer Galanti, Sagie Benaim, and Lior Wolf. Risk bounds for unsupervised cross-domain mapping with ipms. *The Journal of Machine Learning Research*, 22(1):4019–4060, 2021.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research (JMLR)*, 17:2096–2030, 2016.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2014.
- Ishaan Gulrajani and Tatsunori Hashimoto. Identifiability conditions for domain adaptation. In *Proceedings of International Conference on Machine Learning (ICML)*, pp. 7982–7997, 2022.
- 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. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017.
- Junshi Huang, Rogerio S Feris, Qiang Chen, and Shuicheng Yan. Cross-domain image retrieval with a dual attribute-aware ranking network. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp. 1062–1070, 2015.
- Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Proceedings of European Conference on Computer Vision (ECCV)*, pp. 172–189, 2018.
- Aapo Hyvärinen and Petteri Pajunen. Nonlinear independent component analysis: Existence and uniqueness results. *Neural networks*, 12(3):429–439, 1999.
- Aapo Hyvarinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ica using auxiliary variables and generalized contrastive learning. In *Proceedings of International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 859–868. PMLR, 2019.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*, pp. 1125–1134, 2017.
- Chanyong Jung, Gihyun Kwon, and Jong Chul Ye. Exploring patch-wise semantic relation for contrastive learning in image-to-image translation tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 18260–18269, 2022.
- 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.
- Junho Kim, Minjae Kim, Hyeonwoo Kang, and Kwanghee Lee. U-GAT-IT: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2020.
- 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 International Conference on Machine Learning (ICML)*, pp. 1857–1865, 2017.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2015.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 25, 2012.

- Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. Unsupervised machine translation using monolingual corpora only. *arXiv preprint arXiv:1711.00043*, 2017.
- Yann LeCun, Corinna Cortes, and CJ Burges. MNIST handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- Xinyang Li, Jie Hu, Shengchuan Zhang, Xiaopeng Hong, Qixiang Ye, Chenglin Wu, and Rongrong Ji. Attribute guided unpaired image-to-image translation with semi-supervised learning. *arXiv preprint arXiv:1904.12428*, 2019.
- Ming-Yu Liu and Oncel Tuzel. Coupled generative adversarial networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 29, 2016.
- Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, 2017.
- Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. Few-shot unsupervised image-to-image translation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (CVPR)*, pp. 10551–10560, 2019.
- Qi Lyu, Xiao Fu, Weiran Wang, and Songtao Lu. Understanding latent correlation-based multiview learning and self-supervision: An identifiability perspective. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2022.
- Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proceedings of International Conference on Machine Learning (ICML)*, volume 30, pp. 3, 2013.
- Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of International Conference on Computer Vision (ICCV)*, pp. 2794–2802, 2017.
- Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for GANs do actually converge? In *Proceedings of International Conference on Machine Learning (ICML)*, pp. 3481–3490. PMLR, 2018.
- Nikita Moriakov, Jonas Adler, and Jonas Teuwen. Kernel of CycleGAN as a principle homogeneous space. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2020.
- Mostafa Mozafari. Bitmoji faces. <https://www.kaggle.com/datasets/mostafamozafari/bitmoji-faces>, 2020. Accessed on September 20th, 2023.
- Aamir Mustafa and Rafał K Mantiuk. Transformation consistency regularization: A semi-supervised paradigm for image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 599–615, 2020.
- Yingxue Pang, Jianxin Lin, Tao Qin, and Zhibo Chen. Image-to-image translation: Methods and applications. *IEEE Transactions on Multimedia*, 24:3859–3881, 2021.
- Taesung Park, Alexei A Efros, Richard Zhang, and Jun-Yan Zhu. Contrastive learning for unpaired image-to-image translation. In *Proceedings of European Conference on Computer Vision (ECCV)*, pp. 319–345, 2020.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proceedings of International Conference on Machine Learning (ICML)*, pp. 8748–8763, 2021.
- Walter Rudin. *Principles of mathematical analysis*, volume 3. McGraw-hill New York, 1976.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2818–2826, 2016.

- Yaniv Taigman, Adam Polyak, and Lior Wolf. Unsupervised cross-domain image generation. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2017.
- Julius Von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably isolates content from style. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pp. 16451–16467, 2021.
- 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 IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*, pp. 8798–8807, 2018.
- Yaxing Wang, Salman Khan, Abel Gonzalez-Garcia, Joost van de Weijer, and Fahad Shahbaz Khan. Semi-supervised learning for few-shot image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4453–4462, 2020.
- Wayne Wu, Kaidi Cao, Cheng Li, Chen Qian, and Chen Change Loy. TransGaGa: Geometry-aware unsupervised image-to-image translation. In *Proceedings of IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*, pp. 8012–8021, 2019.
- Shaoan Xie, Qirong Ho, and Kun Zhang. Unsupervised image-to-image translation with density changing regularization. *Advances in Neural Information Processing Systems (NeurIPS)*, 35: 28545–28558, 2022.
- Yanwu Xu, Shaoan Xie, Wenhao Wu, Kun Zhang, Mingming Gong, and Kayhan Batmanghelich. Maximum spatial perturbation consistency for unpaired image-to-image translation. In *Proceedings of IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*, pp. 18311–18320, 2022.
- Shuai Yang, Liming Jiang, Ziwei Liu, and Chen Change Loy. Gp-unit: Generative prior for versatile unsupervised image-to-image translation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- Yasin Yaz, Chuan-Sheng Foo, Stefan Winkler, Kim-Hui Yap, Georgios Piliouras, Vijay Chandrasekhar, et al. The unusual effectiveness of averaging in gan training. In *International Conference on Learning Representations*, 2018.
- 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/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 586–595, 2018.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of IEEE/CVF Computer Vision and Pattern Recognition (CVPR)*, pp. 2223–2232, 2017.
- Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1): 43–76, 2020.
- Roland S Zimmermann, Yash Sharma, Steffen Schneider, Matthias Bethge, and Wieland Brendel. Contrastive learning inverts the data generating process. In *Proceedings of International Conference on Machine Learning (ICML)*, pp. 12979–12990, 2021.

Supplementary Material of “Towards Identifiable Unsupervised Domain Translation: A Diversified Distribution Matching Approach”

A PRELIMINARIES

A.1 NOTATION

- x , \mathbf{x} , \mathcal{X} denote a scalar, vector, and a set, respectively.
- $p(\mathbf{x})$ and $p(\mathbf{x}|u)$ denote the marginal *probability density function* (PDF) of \mathbf{x} and conditional PDF of \mathbf{x} conditioned on u , respectively.
- $\|\mathbf{x}\|_2$ denotes the ℓ_2 -norm of \mathbf{x} .
- $\text{dia}(\mathcal{A}) = \sup_{\mathbf{a}, \mathbf{b} \in \mathcal{A}} \|\mathbf{a} - \mathbf{b}\|_2$.
- $\mathbb{I} : \mathcal{X} \rightarrow \mathcal{X}$ denotes the identity function such that $\mathbb{I}(\mathbf{x}) = \mathbf{x}, \forall \mathbf{x} \in \mathcal{X}$.
- \mathcal{A}^c , $\text{cl}(\mathcal{A})$, $\text{bd}(\mathcal{A})$ and $\text{int}(\mathcal{A})$ denote the complement, closure, boundary, and the interior of set \mathcal{A} .
- A set \mathcal{A} is said to have strictly positive measure under $p(\mathbf{x})$ if and only if $\mathbb{P}_{\mathbf{x}}[\mathcal{A}] > 0$.
- For a (random) vector \mathbf{x} , $x(i)$ and $[\mathbf{x}]_i$ denote the i th element of \mathbf{x} , and $\mathbf{x}(i : j)$ denotes $[x(i), x(i+1), \dots, x(j)]$.
- Distance between two sets is defined as

$$\text{dist}(\mathcal{A}, \mathcal{B}) = \inf_{\mathbf{a} \in \mathcal{A}, \mathbf{b} \in \mathcal{B}} \|\mathbf{a} - \mathbf{b}\|_2.$$

- Distance between a set and a point is defined as

$$\text{dist}(\mathbf{a}, \mathcal{B}) = \inf_{\mathbf{b} \in \mathcal{B}} \|\mathbf{a} - \mathbf{b}\|_2.$$

- $\mathcal{N}_\epsilon(\mathbf{z})$ denotes the ϵ -neighborhood of $\mathbf{z} \in \mathbb{R}^N$ defined as

$$\mathcal{N}_\epsilon(\mathbf{z}) = \{\hat{\mathbf{z}} \in \mathbb{R}^N \mid \|\mathbf{z} - \hat{\mathbf{z}}\|_2 < \epsilon\}.$$

- $\text{conn}(\mathcal{A})$ denotes the set of connected components of \mathcal{A} (see definition of connected components in Appendix A.2).
- For any function $\mathbf{m} : \mathcal{W} \rightarrow \mathcal{Z}$, and set $\mathcal{A} \subseteq \mathcal{W}$, $\mathbf{m}(\mathcal{A}) = \{\mathbf{m}(\mathbf{w}) \in \mathcal{Z} \mid \mathbf{w} \in \mathcal{A}\}$

A.2 DEFINITIONS

We will employ standard notions from real analysis. We refer the readers to (Carothers, 2000; Rudin, 1976) for precise definitions and more details. Here we provide working definition with illustration.

Connected set. A set \mathcal{C} is connected (in \mathcal{X}), if and only if there does not exist any disjoint non-empty open sets $\mathcal{A}, \mathcal{B} \subset \mathcal{X}$ such that $\mathcal{A} \cap \mathcal{C} \neq \emptyset, \mathcal{B} \cap \mathcal{C} \neq \emptyset$, and $\mathcal{C} \subset \mathcal{A} \cup \mathcal{B}$ (see Fig. 9).

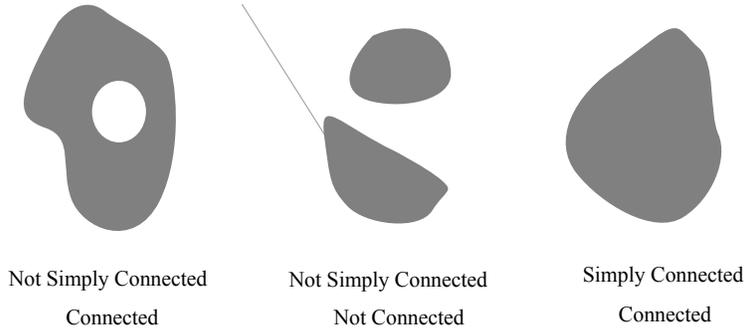


Figure 9: Illustration of connected and simply connected sets

Simply connected set. A simply connected set is a connected set such that any simple closed curve can be shrunk to a point continuously in the set (see Fig. 9).

Connected components. Given a set \mathcal{A} , the maximal connected subsets of \mathcal{A} , such that the subsets are not themselves contained in any other connected subsets of \mathcal{A} , are called connected components of \mathcal{A} . Specifically, a connected set $\mathcal{C} \subseteq \mathcal{A}$ is a connected component of \mathcal{A} if there does not exist any other connected set $\mathcal{D} \subseteq \mathcal{A}$, such that $\mathcal{C} \subset \mathcal{D}$. In Fig. 10, \mathcal{A} denotes the entire shaded regions, and has three connected components \mathcal{C}_1 , \mathcal{C}_2 , and \mathcal{C}_3 . Note that any set can be uniquely written as a disjoint union of its connected components. In Fig. 10, $\mathcal{C}_1 \cup \mathcal{C}_2 \cup \mathcal{C}_3$ is a unique disjoint union representing \mathcal{A} .

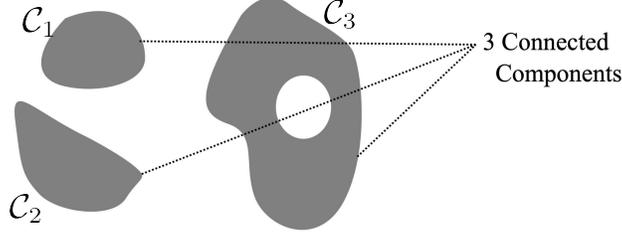


Figure 10: A set $\mathcal{A} = \mathcal{C}_1 \cup \mathcal{C}_2 \cup \mathcal{C}_3$ with 3 connected components: \mathcal{C}_1 , \mathcal{C}_2 , and \mathcal{C}_3 .

Continuous function. A function $m : \mathcal{W} \rightarrow \mathcal{Z}$, with $\mathcal{W} \subseteq \mathbb{R}^W$, $\mathcal{Z} \subseteq \mathbb{R}^Z$ is said to be continuous if for any $w \in \mathcal{W}$ and $\epsilon > 0$, there exists a $\delta > 0$ such that

$$m((\mathcal{N}_\delta(w) \cap \mathcal{W})) \subset \mathcal{N}_\epsilon(m(w)) \cap \mathcal{Z}.$$

Continuous and invertible Functions. If a function $m : \mathcal{W} \rightarrow \mathcal{Z}$ is continuous and invertible, then its inverse $m^{-1} : \mathcal{Z} \rightarrow \mathcal{W}$ is also continuous. Some useful properties of continuous and invertible function m are as follows:

- If $\mathcal{A} \subseteq \mathcal{W}$ is closed, then $m(\mathcal{A})$ is also closed.
- If $\mathcal{A} \subseteq \mathcal{W}$ is open, then $m(\mathcal{A})$ is also open.
- If $\mathcal{A} \subseteq \mathcal{W}$ is connected, then $m(\mathcal{A})$ is also connected.

B PROOF OF LEMMAS AND FACTS

Note that for the ease of reading, the lemmas, facts, and theorems from the main paper are re-stated and highlighted using shaded boxes.

Proposition 1. *Suppose that \mathbb{P}_x admits a continuous PDF, $p(x)$ and $p(x) > 0, \forall x \in \mathcal{X}$. Assume that \mathcal{X} is simply connected. Then, there exists a continuous non-trivial (non-identity) $h(\cdot)$ such that $h_{\#\mathbb{P}_x} = \mathbb{P}_x$.*

Proof. We want to show that there exists a continuous $h : \mathcal{X} \rightarrow \mathcal{X}$ such that

$$h_{\#\mathbb{P}_x} = \mathbb{P}_x.$$

To this end, we will construct such MPA by reducing the problem of finding an MPA of $p(x)$ to finding an MPA of the uniform distribution. Note that one can always construct a continuous invertible function $d : \mathcal{X} \rightarrow (0, 1)^{D_x}$, such that the function maps any continuous distribution with a simply connected support to the uniform distribution. This mapping can be found via the so-called *Darmois construction* (Darmois, 1951; Hyvärinen & Pajunen, 1999). Specifically, under the Darmois construction, the i th output of $d(\bar{x})$, $\forall \bar{x} \in \mathcal{X}$ is given by

$$[d(\bar{x})]_i := F(\bar{x}(i) \mid \mathbf{x}(1:i-1) = \bar{\mathbf{x}}(1:i-1)), \quad i = 1, \dots, N,$$

where $F(\bar{x}(i) \mid \cdot)$ denotes the conditional CDF of $x(i)$, i.e,

$$F(\bar{x}(i) \mid \mathbf{x}(1:i-1) = \bar{\mathbf{x}}(1:i-1)) = \mathbb{P}_{x(i) \mid \mathbf{x}(1:i-1) = \bar{\mathbf{x}}(1:i-1)} [\{x(i) : x(i) \leq \bar{x}(i)\}];$$

see more detailed introduction to the Darmois construction in (Hyvärinen & Pajunen, 1999).

With the constructed \mathbf{d} , one can form a continuous mapping $\mathbf{h} : \mathcal{X} \rightarrow \mathcal{X}$ as follows

$$\mathbf{h} = \mathbf{d}^{-1} \circ \mathbf{h}_U \circ \mathbf{d},$$

where $\mathbf{h}_U : (0, 1)^{D_x} \rightarrow (0, 1)^{D_x}$ is a continuous MPA on the uniform distribution over $(0, 1)^{D_x}$. Since \mathbf{d} is continuous for a continuous distribution, \mathbf{h} is continuous because it is the composition of continuous functions.

Now, it remains to show that \mathbf{h}_U exists. A simple example of \mathbf{h}_U is reflection around the mean of $d(\mathbf{x})$, i.e.,

$$\mathbf{h}_U(\mathbf{z}) = -\mathbf{z} + 2\boldsymbol{\mu},$$

where $\boldsymbol{\mu} = [1/2, \dots, 1/2]^\top \in \mathbb{R}^{D_x}$. This concludes the proof. \square

Lemma 1. Assume that an optimal solution of (7) is $(\hat{\mathbf{f}}, \hat{\mathbf{g}}, \{\hat{\mathbf{d}}_x^{(i)}, \hat{\mathbf{d}}_y^{(i)}\})$. Then, under Assumption 1, we have $\mathbb{P}_{\mathbf{x}|u=u_i} = \hat{\mathbf{f}}_{\#\mathbb{P}_{\mathbf{y}|u=u_i}}$, $\mathbb{P}_{\mathbf{y}|u=u_i} = \hat{\mathbf{g}}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}}$, $\forall i \in [I]$, and $\hat{\mathbf{f}} = \hat{\mathbf{g}}^{-1}$, a.e.

Proof. Fact 1 is a direct consequence of (Goodfellow et al., 2014, Theorem 1).

First of all, recall the objective in (7):

$$\min_{\mathbf{f}, \mathbf{g}} \max_{\{\mathbf{d}_x^{(i)}, \mathbf{d}_y^{(i)}\}} \sum_{i=1}^I \left(\mathcal{L}_{\text{DSGAN}}(\mathbf{g}, \mathbf{d}_y^{(i)}, \mathbf{x}, \mathbf{y}) + \mathcal{L}_{\text{DSGAN}}(\mathbf{f}, \mathbf{d}_x^{(i)}, \mathbf{x}, \mathbf{y}) \right) + \lambda \mathcal{L}_{\text{cyc}}(\mathbf{g}, \mathbf{f}). \quad (8)$$

The global minimum of $\mathcal{L}_{\text{DSGAN}}(\mathbf{g}, \mathbf{d}_y^{(i)}, \mathbf{x}, \mathbf{y})$ is achieved when (Goodfellow et al., 2014, Theorem 1)

$$\mathbf{g}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} = \mathbb{P}_{\mathbf{y}|u=u_i}.$$

Similarly, the global minimum of $\mathcal{L}_{\text{DSGAN}}(\mathbf{f}, \mathbf{d}_x^{(i)}, \mathbf{x}, \mathbf{y})$ is achieved when

$$\mathbf{f}_{\#\mathbb{P}_{\mathbf{y}|u=u_i}} = \mathbb{P}_{\mathbf{x}|u=u_i}.$$

Finally, the global minimum of $\mathcal{L}_{\text{cyc}}(\mathbf{g}, \mathbf{f})$, which is zero, is achieved when

$$\mathbf{g} = \mathbf{f}^{-1}, \text{ a.e.}$$

We know that \mathbf{g}^* and \mathbf{f}^* can achieve global minimums of all loss terms simultaneously. Hence the solution of (7), $\hat{\mathbf{f}}$ and $\hat{\mathbf{g}}$, should satisfy

$$\mathbb{P}_{\mathbf{x}|u=u_i} = \hat{\mathbf{f}}_{\#\mathbb{P}_{\mathbf{y}|u=u_i}}, \mathbb{P}_{\mathbf{y}|u=u_i} = \hat{\mathbf{g}}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}}, \forall i \in [I], \text{ and } \hat{\mathbf{f}} = \hat{\mathbf{g}}^{-1}, \text{ a.e.}$$

\square

Fact 2. Suppose that Assumption 1 holds. Then, there exists an admissible MPA of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ if and only if there exists an admissible MPA of $\{\mathbb{P}_{\mathbf{y}|u=u_i}\}_{i=1}^I$.

Proof. Let \mathbf{h} be an admissible MPA of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$. Then

$$\begin{aligned} \mathbf{h}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} &= \mathbb{P}_{\mathbf{x}|u=u_i}, \forall i \in [I]. \\ \iff \mathbf{g}^* \circ \mathbf{h}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} &= \mathbf{g}^*_{\#\mathbb{P}_{\mathbf{x}|u=u_i}}, \forall i \in [I]. \\ \iff \mathbf{g}^* \circ \mathbf{h} \circ \mathbf{f}^*_{\#\mathbb{P}_{\mathbf{y}|u=u_i}} &= \mathbb{P}_{\mathbf{y}|u=u_i}, \forall i \in [I]. \end{aligned}$$

This implies that $\mathbf{g}^* \circ \mathbf{h} \circ \mathbf{f}^*$ is an admissible MPA of $\{\mathbb{P}_{\mathbf{y}|u=u_i}\}_{i=1}^I$ if and only if \mathbf{h} is an admissible MPA of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$.

Hence, there exists an admissible MPA of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ if and only if there exists an admissible MPA of $\{\mathbb{P}_{\mathbf{y}|u=u_i}\}_{i=1}^I$. \square

C PROOF OF THEOREMS

C.1 PROOF OF THEOREM 1

Theorem 1. *Suppose that Assumption 1 holds. Let $E_{i,j}$ denote the event that the set $\{\mathbb{P}_{\mathbf{x}|u=u_i}, \mathbb{P}_{\mathbf{x}|u=u_j}\}$ does not satisfy the SDC. Assume that $\Pr[E_{i,j}] \leq \rho$ for any $i \neq j$, where $i, j \in [I]$. Let $(\hat{\mathbf{f}}, \hat{\mathbf{g}})$ be from an optimal solution of the DIMENSION loss (7). Then, there is no admissible MPA of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ of the solution, i.e., $\hat{\mathbf{f}} = \mathbf{f}^*$, a.e., $\hat{\mathbf{g}} = \mathbf{g}^*$, a.e., with a probability of at least $1 - \rho^{\binom{I}{2}}$.*

Theorem 1 is a direct consequence of following lemma:

Lemma A.1. *Suppose that Assumption 1 holds. Assume that $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ are sufficiently diverse. Then, $\hat{\mathbf{g}} = \mathbf{g}^*$ and $\hat{\mathbf{f}} = \mathbf{f}^*$, a.e.*

Proof of Lemma A.1. First, we show that no non-trivial continuous admissible MPA exists for $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$, i.e., if a continuous \mathbf{h} satisfies

$$\mathbf{h}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} = \mathbb{P}_{\mathbf{x}|u=u_i} \forall i \in [I], \quad (9)$$

then $\mathbf{h} = \mathbb{1}$, a.e.

Eq. (9), by the definition of push-forward measure, implies that

$$\implies \mathbb{P}_{\mathbf{x}|u_i}[\mathbf{h}(\mathcal{A})] = \mathbb{P}_{\mathbf{x}|u_i}[\mathcal{A}], \forall i \in [I]. \quad (10)$$

For the sake of contradiction assume that \mathbf{h} satisfies (9), however, $\mathbf{h} \neq \mathbb{1}$ on a set of strictly positive measure. This means that there exists a $\bar{\mathbf{x}} \in \mathcal{X}$ such that

$$\mathbf{h}(\bar{\mathbf{x}}) \neq \bar{\mathbf{x}}.$$

Now, let us define an open set around $\bar{\mathbf{x}}$ denoted by \mathcal{D} such that

$$\mathcal{D} = \mathcal{N}_d(\bar{\mathbf{x}}) \cap \mathcal{X}.$$

Because of the continuity and invertibility of \mathbf{h} , $\mathbf{h}(\mathcal{D}) \subseteq \mathcal{X}$ is also an open set and

$$\mathbf{h}(\bar{\mathbf{x}}) \in \mathbf{h}(\mathcal{D}).$$

Now, one can select d to be small enough (because of the continuity of \mathbf{h}) such that $\mathcal{D} \cap \mathbf{h}(\mathcal{D}) = \emptyset$ and \mathcal{D} is a connected set. \mathcal{D} being a connected set implies that $\mathbf{h}(\mathcal{D})$ is also connected.

The above is a contradiction to Assumption 4 since \mathcal{D} and $\mathbf{h}(\mathcal{D})$ are two disjoint, open and connected sets which satisfy

$$\mathbb{P}_{\mathbf{x}|u_i}[\mathbf{h}(\mathcal{D})] = \mathbb{P}_{\mathbf{x}|u_i}[\mathcal{D}], \forall i \in [I].$$

Hence, any \mathbf{h} that satisfy $\mathbf{h}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} = \mathbb{P}_{\mathbf{x}|u=u_i}$ is such that $\mathbf{h} = \mathbb{1}$, a.e.

Finally, We want to show that $\hat{\mathbf{g}} = \mathbf{g}^*$, a.e. Lemma 1 implies that

$$\begin{aligned} \hat{\mathbf{g}}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} &= \mathbb{P}_{\mathbf{y}|u=u_i}, \forall i \in [I] \\ \implies \hat{\mathbf{g}}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} &= \mathbf{g}^*_{\#\mathbb{P}_{\mathbf{x}|u=u_i}}, \forall i \in [I] \\ \stackrel{(a)}{\implies} \mathbf{g}^{*-1} \circ \hat{\mathbf{g}}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} &= \mathbb{P}_{\mathbf{x}|u=u_i}, \forall i \in [I] \end{aligned} \quad (11)$$

where (a) is obtained by applying \mathbf{g}^{*-1} on both sides, which is allowed because applying the same function preserves the equivalence of the distributions.

Eq. (11) implies that $\mathbf{g}^* \circ \hat{\mathbf{g}}$ is a continuous admissible MPA of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$, which means that the following has to hold:

$$\mathbf{g}^{*-1} \circ \hat{\mathbf{g}} = \mathbb{1}, \text{ a.e.}$$

Therefore, we always have $\hat{\mathbf{g}} = \mathbf{g}^*$, a.e. By role symmetry of \mathbf{f}^* and \mathbf{g}^* (also see Fact 2), we also have $\hat{\mathbf{f}} = \mathbf{f}^*$, a.e. \square

Proof of Theorem 1. Using the assumption that $\Pr[\mathbb{E}_{i,j}] \leq \rho$, the probability that $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ are not sufficiently diverse can be bounded as follows:

$$\begin{aligned} & \Pr[\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I \text{ are not sufficiently diverse}] \\ & \stackrel{(a)}{\leq} \Pr \left[\bigcap_{i,j \in [I], i < j} \mathbb{E}_{i,j} \right] \\ & \stackrel{(b)}{=} \bigcap_{i,j \in [I], i < j} \Pr[\mathbb{E}_{i,j}] \\ & \leq \rho^{\binom{I}{2}}, \end{aligned}$$

where the (a) holds since $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ not being sufficiently diverse implies the existence of open connected sets \mathcal{A} and \mathcal{B} such that

$$\begin{aligned} & \mathbb{P}_{\mathbf{x}|u_i}[\mathcal{A}] = \mathbb{P}_{\mathbf{x}|u_i}[\mathcal{B}], \forall i \in [I] \\ \implies & \mathbb{P}_{\mathbf{x}|u_i}[\mathcal{A}] = \mathbb{P}_{\mathbf{x}|u_i}[\mathcal{B}], \forall i \in \{i, j\} \subset [I]. \end{aligned}$$

Finally, (b) is due to the independence of the events $\mathbb{E}_{i,j}$ and $\mathbb{E}_{i,j'}$ for $j \neq j'$.

Hence, $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ are sufficiently diverse with probability at least $1 - \rho^{\binom{I}{2}}$, which implies that $\hat{\mathbf{f}} = \mathbf{f}^*$ and $\hat{\mathbf{g}} = \mathbf{g}^*$ with probability at least $1 - \rho^{\binom{I}{2}}$. \square

C.2 PROOF OF THEOREM 2

Theorem 2 (Robust Identifiability). *Suppose that Assumption 1 holds with \mathbf{g}^* being L -Lipschitz continuous, and that any pair of $(\mathbb{P}_{\mathbf{x}|u_i}, \mathbb{P}_{\mathbf{x}|u_j})$ satisfies the r -SDC (cf. Definition 5) with probability at least $1 - \gamma$, i.e., $\Pr[M_{i,j} \geq r] \leq \gamma$ for any $i \neq j$, where $(i, j) \in [I] \times [J]$. Let $\hat{\mathbf{g}}$ be from any optimal solution of the DIMENSION loss in (7). Then, we have $\|\hat{\mathbf{g}}(\mathbf{x}) - \mathbf{g}^*(\mathbf{x})\|_2 \leq 2rL$, $\forall \mathbf{x} \in \mathcal{X}$, with a probability of at least $1 - \gamma^{\binom{I}{2}}$. The same holds for $\hat{\mathbf{f}}$.*

First, consider the following lemma.

Lemma A.2. *Given any continuous admissible MPA \mathbf{h} of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$, let $\mathcal{E}_{\mathbf{h}}$ be a set defined as*

$$\mathcal{E}_{\mathbf{h}} = \{\mathbf{x} \mid \mathbf{h}(\mathbf{x}) \neq \mathbf{x}, \forall \mathbf{x} \in \mathcal{X}\}.$$

Then, any connected component $\mathcal{C} \subseteq \text{cl}(\mathcal{E}_{\mathbf{h}})$ satisfies

$$\mathbf{x} \in \mathcal{C} \implies \mathbf{h}(\mathbf{x}) \in \mathcal{C}.$$

Lemma A.2 states an interesting property of a subset of \mathcal{X} (namely, $\mathcal{E}_{\mathbf{h}}$) that is “modified” by the continuous MPA \mathbf{h} . Here, “modification” means that any point in the subset will land on a different point after the \mathbf{h} -transformation. The lemma shows that the source point from $\mathcal{E}_{\mathbf{h}}$ and its \mathbf{h} -transformation both reside in the same connected component, namely, \mathcal{C} . This will be useful in proving Theorem 2.

Proof Idea: The main idea behind the proof of Lemma A.2 is to first note that any point outside of $\text{cl}(\mathcal{E}_{\mathbf{h}})$ is stationary under the transformation \mathbf{h} (i.e., $\mathbf{h}(\mathbf{x}) = \mathbf{x}$). Next, if there was a point $\bar{\mathbf{x}}$ from a connected component $\mathcal{C}_1 \in \text{conn}(\text{cl}(\mathcal{E}_{\mathbf{h}}))$ was such that $\mathbf{h}(\bar{\mathbf{x}})$ was not in \mathcal{C}_1 , then it should be either in $\mathcal{X} \setminus \text{cl}(\mathcal{E}_{\mathbf{h}})$ or in $\text{cl}(\mathcal{E}_{\mathbf{h}}) \setminus \mathcal{C}_1$. However, since \mathbf{h} is invertible, \mathbf{h} cannot map a point from \mathcal{C}_1 to a $\mathcal{X} \setminus \text{cl}(\mathcal{E}_{\mathbf{h}})$. Therefore $\mathbf{h}(\bar{\mathbf{x}})$ should lie in $\mathcal{E}_{\mathbf{h}} \setminus \mathcal{C}_1$. But this will make the function \mathbf{h} discontinuous. Hence, $\mathbf{h}(\bar{\mathbf{x}})$ should be in \mathcal{C}_1 .

Proof of Lemma A.2. Let $\text{conn}(\text{cl}(\mathcal{E}_{\mathbf{h}}))$ denote the set of connected components of $\text{cl}(\mathcal{E}_{\mathbf{h}})$. Suppose that there exists $\bar{\mathbf{x}} \in \mathcal{C}_1$ and $\mathcal{C}_1 \in \text{conn}(\text{cl}(\mathcal{E}_{\mathbf{h}}))$ such that $\mathbf{h}(\bar{\mathbf{x}}) \notin \mathcal{C}_1$. First,

$$\begin{aligned} \mathbf{h}(\tilde{\mathbf{x}}) = \tilde{\mathbf{x}}, \forall \tilde{\mathbf{x}} \in \mathcal{X} \setminus \text{cl}(\mathcal{E}_{\mathbf{h}}) & \stackrel{(a)}{\implies} \mathbf{h}(\bar{\mathbf{x}}) \neq \tilde{\mathbf{x}}, \forall \tilde{\mathbf{x}} \in \mathcal{X} \setminus \text{cl}(\mathcal{E}_{\mathbf{h}}) \\ & \implies \mathbf{h}(\bar{\mathbf{x}}) \in \mathcal{C}_2, \text{ for some } \mathcal{C}_2 \in \text{conn}(\text{cl}(\mathcal{E}_{\mathbf{h}})) \setminus \mathcal{C}_1, \end{aligned}$$

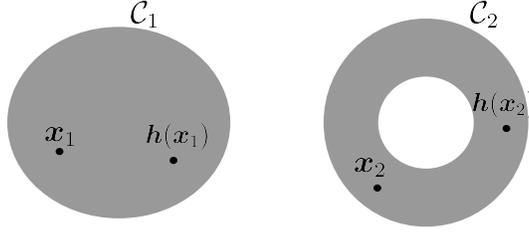


Figure 11: Illustration of Lemma A.2. $\mathcal{C}_1, \mathcal{C}_2$ are the two connected components of $\text{cl}(\mathcal{E}_h)$. In this case, $\text{cl}(\mathcal{E}_h) = \mathcal{C}_1 \cup \mathcal{C}_2$. Points x_1 and x_2 inside \mathcal{C}_1 and \mathcal{C}_2 stay inside the same connected component after transformation by h .

where (a) is due to the invertibility of h .

Because of the continuity of h , the set $h(\mathcal{C}_1)$ is a closed connected set containing $h(\bar{x})$. However, $h(\mathcal{C}_1) \cap (\mathcal{X} \setminus \text{cl}(\mathcal{E}_h)) = \emptyset$. This means that

$$h(\mathcal{C}_1) \subseteq \mathcal{C}_2, \quad (12)$$

otherwise $h(\mathcal{C}_1)$ would be disconnected.

Note that \mathcal{C}_1 and \mathcal{C}_2 are closed, connected and disjoint sets (by the property of connected components). Therefore, one can define ϵ as follows:

$$\epsilon := \text{dist}(\mathcal{C}_1, \mathcal{C}_2) > 0. \quad (13)$$

Now, take any point $x_b \in \text{bd}(\mathcal{C}_1)$, where $\text{bd}(\mathcal{C}_1)$ denotes the boundary of \mathcal{C}_1 . Due to the continuity of h , there exists a $\delta > 0$ such that

$$h(\mathcal{N}_\delta(x_b)) \subseteq \mathcal{N}_{\epsilon/4}(h(x_b)). \quad (14)$$

However, take any point $z \in \mathcal{N}_\delta(x_b) \setminus \text{cl}(\mathcal{E}_h)$ with $\|z - x_b\|_2 < \epsilon/4$. Such a point exists because any neighborhood of a point on the boundary of a closed set has a non-empty intersection with the complement of the closed set. Therefore, we have

$$h(z) = z \text{ because } z \notin \mathcal{E}_h. \quad (15)$$

Since (12) implies that $h(x_b) \in \mathcal{C}_2$,

$$\|h(x_b) - x_b\|_2 \geq \epsilon \quad \text{and} \quad \text{dist}(x_b, \mathcal{N}_{\epsilon/4}(h(x_b))) \geq \frac{3\epsilon}{4}.$$

Therefore

$$\begin{aligned} \text{dist}(x_b, \mathcal{N}_{\epsilon/4}(h(x_b))) &\leq \|x_b - z\|_2 + \text{dist}(z, \mathcal{N}_{\epsilon/4}(h(x_b))) \\ &\implies \frac{3\epsilon}{4} \leq \frac{\epsilon}{4} + \text{dist}(z, \mathcal{N}_{\epsilon/4}(h(x_b))) \\ &\implies \text{dist}(z, \mathcal{N}_{\epsilon/4}(h(x_b))) \geq \frac{\epsilon}{2} \\ &\implies \text{dist}(h(z), \mathcal{N}_{\epsilon/4}(h(x_b))) \geq \frac{\epsilon}{2}, \end{aligned} \quad (16)$$

where (16) is by (15). Note that (16) is a contradiction to (14). Hence, we have

$$x \in \mathcal{C} \text{ for any } \mathcal{C} \in \text{conn}(\text{cl}(\mathcal{E}_h)) \implies h(x) \in \mathcal{C}.$$

This concludes the proof. \square

Lemma A.3. *Let g^* be L -Lipschitz continuous. Suppose that any pair of $(\mathbb{P}_{x|u_i}, \mathbb{P}_{x|u_j})$ satisfies the r -SDC (cf. Definition 5). Then*

$$\|\hat{g}(x) - g^*(x)\|_2 \leq 2rL. \quad (17)$$

Proof Idea: The proof is by contradiction. Suppose that under the conditions of Lemma A.3, Eq. (17) does not hold for some $\bar{\mathbf{x}} \in \mathcal{X}$. Then, there would exist a continuous non-trivial admissible MPA $\bar{\mathbf{h}}$ of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$ such that $\|\bar{\mathbf{h}}(\bar{\mathbf{x}}) - \bar{\mathbf{x}}\|_2 > 2r$. However, this would imply, using Lemma A.2, that one can construct an open, connected, disjoint set pair $(\mathcal{A}, \mathcal{B})$ whose diameters are large, which is a contradiction to r -SDC that $M \leq r$.

Proof of Lemma A.3. Suppose that there exists $\bar{\mathbf{x}} \in \mathcal{X}$ such that

$$\|\hat{\mathbf{g}}(\bar{\mathbf{x}}) - \mathbf{g}^*(\bar{\mathbf{x}})\|_2 > 2rL. \quad (18)$$

Eq. (18) means that $\hat{\mathbf{g}} \neq \mathbf{g}^*$. By Lemma 1, we have that

$$\begin{aligned} \hat{\mathbf{g}}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} &= \mathbb{P}_{\mathbf{y}|u=u_i} \\ \iff \hat{\mathbf{g}}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} &= \mathbf{g}^*_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} \\ \iff \mathbf{g}^{*-1} \circ \hat{\mathbf{g}}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} &= \mathbf{g}^{*-1} \circ \mathbf{g}^*_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} \\ \iff \mathbf{g}^{*-1} \circ \hat{\mathbf{g}}_{\#\mathbb{P}_{\mathbf{x}|u=u_i}} &= \mathbb{P}_{\mathbf{x}|u=u_i} \end{aligned}$$

As $\hat{\mathbf{g}} \neq \mathbf{g}^*$, the function $\bar{\mathbf{h}} := \mathbf{g}^{*-1} \circ \hat{\mathbf{g}} \neq \mathbb{I}$ is a continuous admissible MPA of $\{\mathbb{P}_{\mathbf{x}|u=u_i}\}_{i=1}^I$. This implies that

$$\hat{\mathbf{g}} = \mathbf{g}^* \circ \bar{\mathbf{h}}. \quad (19)$$

Using (19), Eq. (18) implies that

$$\begin{aligned} \|\mathbf{g}^* \circ \bar{\mathbf{h}}(\bar{\mathbf{x}}) - \mathbf{g}^*(\bar{\mathbf{x}})\|_2 &> 2rL \\ \implies L\|\bar{\mathbf{h}}(\bar{\mathbf{x}}) - \bar{\mathbf{x}}\|_2 &> 2rL \\ \implies \|\bar{\mathbf{h}}(\bar{\mathbf{x}}) - \bar{\mathbf{x}}\|_2 &> 2r. \end{aligned} \quad (20)$$

Note that one can re-express (20) as

$$\|\bar{\mathbf{h}}(\bar{\mathbf{x}}) - \bar{\mathbf{x}}\|_2 = 2r + \epsilon, \quad (21)$$

using a certain $\epsilon > 0$. By Lemma A.2, we know that

$$\bar{\mathbf{x}} \in \bar{\mathcal{C}}, \quad \bar{\mathbf{h}}(\bar{\mathbf{x}}) \in \bar{\mathcal{C}},$$

where $\bar{\mathcal{C}}$ is a connected component of $\text{cl}(\mathcal{E}_{\bar{\mathbf{h}}})$.

Now, let $d > 0$ and

$$\mathcal{T}_d = \mathcal{N}_d(\bar{\mathbf{x}}) \cap \bar{\mathcal{C}}.$$

Let \mathcal{R}_d denote the connected component of \mathcal{T}_d that contains $\bar{\mathbf{x}}$. Note that we need to consider the connected component as \mathcal{T}_d can be a disconnected set. An illustration of these sets can be seen in Fig. 12.

From (21), it is easy to see that $\text{dia}(\bar{\mathcal{C}}) \geq 2r + \epsilon$. Then, when $d \leq 2r + \epsilon$, since $\bar{\mathcal{C}}$ is connected, $\text{bd}(\mathcal{N}_d(\bar{\mathbf{x}})) \cap \bar{\mathcal{C}} \neq \emptyset$. Note that the connected property $\bar{\mathcal{C}}$ is necessary for $\text{bd}(\mathcal{N}_d(\bar{\mathbf{x}})) \cap \bar{\mathcal{C}} \neq \emptyset$ to hold for any $d \leq 2r + \epsilon$. One can further select $\mathbf{w} \in \text{bd}(\mathcal{N}_d(\bar{\mathbf{x}})) \cap \bar{\mathcal{C}}$ such that $\mathbf{w} \in \text{cl}(\mathcal{R}_d)$, i.e., \mathbf{w} lies in the same connected component of \mathcal{T}_d as $\bar{\mathbf{x}}$.

Note that such a \mathbf{w} has to exist. Suppose that such a \mathbf{w} does not exist. Then, $\text{cl}(\mathcal{R}_d) \cap \text{bd}(\mathcal{N}_d) = \emptyset$, which means that \mathcal{R}_d would be disconnected from $\bar{\mathcal{C}} \setminus \mathcal{N}_d$. By the definition of \mathcal{R}_d , \mathcal{R}_d is then disconnected from $\mathcal{T}_d \setminus \mathcal{R}_d$, which implies that $\bar{\mathcal{C}}$ —that is a union of \mathcal{R}_d , $\mathcal{T}_d \setminus \mathcal{R}_d$, and $\bar{\mathcal{C}} \setminus \mathcal{N}_d$ —is disconnected. This is a contradiction. Hence $\text{cl}(\mathcal{R}_d) \cap \text{bd}(\mathcal{N}_d) \neq \emptyset$ holds.

One can see that

$$\|\mathbf{w} - \bar{\mathbf{x}}\|_2 = d, \text{ and } \mathbf{w}, \bar{\mathbf{x}} \in \text{cl}(\mathcal{R}_d) \implies \text{dia}(\mathcal{R}_d) \geq d, \text{ for } d \leq 2r + \epsilon.$$

Hence, there exists a large enough $d \leq 2r + \epsilon$ such that

$$0 < \text{dist}(\mathcal{R}_d, \bar{\mathbf{h}}(\mathcal{R}_d)) < \epsilon/3.$$

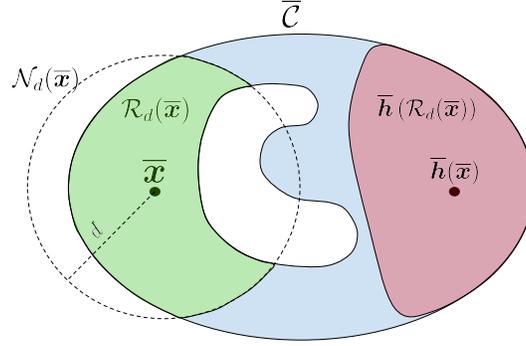


Figure 12: Illustration of the idea in the proof of Lemma A.3. The green shaded region denote \mathcal{R}_d . Note that $\mathcal{T}_d = \mathcal{N}_d(\bar{\mathbf{x}}) \cap \bar{\mathcal{C}}$ is disconnected in this case.

This implies that

$$\max\{\text{dia}(\mathcal{R}_d), \text{dia}(\bar{\mathbf{h}}(\mathcal{R}_d))\} \geq r + \epsilon/3. \quad (22)$$

Indeed, suppose that $\max\{\text{dia}(\mathcal{R}_d), \text{dia}(\bar{\mathbf{h}}(\mathcal{R}_d))\} < r + \epsilon/3$. Then, since $\bar{\mathbf{x}} \in \mathcal{R}_d$ and $\bar{\mathbf{h}}(\bar{\mathbf{x}}) \in \bar{\mathbf{h}}(\mathcal{R}_d)$,

$$\begin{aligned} \|\bar{\mathbf{x}} - \bar{\mathbf{h}}(\bar{\mathbf{x}})\|_2 &\leq 2\max\{\text{dia}(\mathcal{R}_d), \text{dia}(\bar{\mathbf{h}}(\mathcal{R}_d))\} + \text{dist}(\mathcal{R}_d, \bar{\mathbf{h}}(\mathcal{R}_d)) \\ 2r + \epsilon &< 2r + 2\epsilon/3 + \epsilon/3 \\ 2r + \epsilon &< 2r + \epsilon, \end{aligned}$$

which is a contradiction. Hence,

$$\max\{\text{dia}(\mathcal{R}_d), \text{dia}(\bar{\mathbf{h}}(\mathcal{R}_d))\} \geq r + \epsilon/3. \quad (23)$$

Fig. 12 provides a simple illustration of the sets. It follow from the continuity and invertibility of $\bar{\mathbf{h}}$ that $\text{dia}(\mathcal{R}_d) = \text{dia}(\text{int}(\mathcal{R}_d))$ and $\text{dia}(\bar{\mathbf{h}}(\mathcal{R}_d)) = \text{dia}(\text{int}(\bar{\mathbf{h}}(\mathcal{R}_d)))$.

By the same argument of reaching (10), $\{\text{int}(\mathcal{R}_d), \bar{\mathbf{h}}(\text{int}(\mathcal{R}_d))\}$ forms a pair of open, connected, disjoint sets such that

$$\mathbb{P}_{\mathbf{x}|u_i}[\text{int}(\mathcal{R}_d)] = \mathbb{P}_{\mathbf{x}|u_i}[\bar{\mathbf{h}}(\text{int}(\mathcal{R}_d))]. \quad (24)$$

Note that (24) and (23) constitute a contradiction to the assumption that

$$M = \max_{(\mathcal{A}, \mathcal{B}) \in \mathcal{V}} \max\{\text{dia}(\mathcal{A}), \text{dia}(\mathcal{B})\} \leq r$$

for any open, connected, and disjoint sets \mathcal{A} and \mathcal{B}^2

Hence, we must have

$$\|\hat{\mathbf{g}}(\mathbf{x}) - \mathbf{g}^*(\mathbf{x})\|_2 \leq 2rL, \forall \mathbf{x} \in \mathcal{X}.$$

This concludes the proof. \square

Proof of Theorem 2. We can bound the probability with which (17) does not hold as follows:

$$\begin{aligned} &\Pr[(17) \text{ does not hold}] \\ &= \Pr[M > r] \\ &\stackrel{(a)}{\leq} \Pr\left[\bigcap_{i,j \in [I], i < j} M_{i,j}\right] \\ &\stackrel{(b)}{=} \bigcap_{i,j \in [I], i < j} \Pr[M_{i,j}] \\ &\leq \gamma^{\binom{I}{2}}, \end{aligned}$$

²Note that such requirements are used to ensure that the sets have nonzero measure. The statements can be simplified by replacing the ‘‘open and non-empty’’ sets in Definition 4 and Assumption 5 with ‘‘measurable sets with positive measures’’.

where (a) follows because $M > r$ implies that $M_{i,j} > r, \forall i \neq j$, and $i, j \in [I]$ holds, and (b) follows from the independence of the events $M_{i,j} > r$.

Hence, with probability at least $1 - \gamma^{\binom{I}{2}}$,

$$\|\widehat{\mathbf{g}}(\mathbf{x}) - \mathbf{g}^*(\mathbf{x})\|_2 \leq 2rL, \forall \mathbf{x} \in \mathcal{X}.$$

The same result follows for $\widehat{\mathbf{f}}$ if \mathbf{f}^* is L -Lipschitz continuous, following the same procedure as above. \square

D ADDITIONAL REMARK: RELATION TO SUPERVISED DOMAIN TRANSLATION

A remark on objective (7) is that supervised domain translation can be seen as a special case of (7). When paired samples $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$ are available, one can view the auxiliary information u as the sample identity. Specifically, $\mathbb{P}_{\mathbf{x}|u=u_i}$ and $\mathbb{P}_{\mathbf{y}|u=u_i}$ are Dirac delta distributions peaked at \mathbf{x}_i and \mathbf{y}_i , respectively. Matching distributions between $\mathbb{P}_{\mathbf{x}|u=u_i}$ and $\mathbf{f}_{\# \mathbb{P}_{\mathbf{y}|u=u_i}}$ will be equivalent to enforcing $\mathbf{x}_i = \mathbf{f}(\mathbf{y}_i)$. Therefore, the sample loss will be equivalent to minimizing the following objective:

$$\underset{\mathbf{f}, \mathbf{g}}{\text{minimize}} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{f}(\mathbf{y}_i)\|_2^2 + \|\mathbf{y}_i - \mathbf{g}(\mathbf{x}_i)\|_2^2,$$

which is exactly the supervised learning loss. This makes the distribution matching problem boil down to a sample matching problem.

E SYNTHETIC DATA EXPERIMENTS

In this section, we use controlled generation to validate our identifiability theorems.

Data Generation. We generate \mathbf{x} from a Gaussian mixture with Q components. Let $\{\mathbb{P}_{\mathbf{x}}^{(q)}\}_{q=1}^Q$ denote the Q component distributions of the Gaussian mixture, i.e.,

$$\mathbb{P}_{\mathbf{x}}^{(q)} \sim \mathcal{N}(\boldsymbol{\mu}_q, \boldsymbol{\Sigma}), \quad q = 1, \dots, Q.$$

Here, each $\boldsymbol{\mu}_q$ is sampled randomly from the uniform distribution in \mathbb{R}^2 , i.e., $\text{Unif}([-1, 1]^2)$. We set the covariance to be $\boldsymbol{\Sigma} = 0.3^2 \mathbf{Q}$. To represent \mathbf{g}^* , we use a three-layer multi-layer perceptron (MLP) with smoothed leaky ReLU, which is defined as $s(x) = \alpha x + (1 - \alpha) \log(1 + \exp(x))$, where we set α to 0.2. To make \mathbf{g}^* invertible, we generate the neural network weights using the same process as in (Hyvärinen & Pajunen, 1999; Zimmermann et al., 2021). Specifically, we use two-hidden units in each layer. We first generate 10,000 2×2 matrices, whose elements are sampled randomly from uniform distribution $\text{Unif}([-1, 1])$. The matrices' columns are normalized by their respective ℓ_2 norms. In addition, only the top 25% well-conditioned matrices in terms of the condition number are used. This way, all the layers of the \mathbf{g}^* are relatively well-conditioned invertible matrices. Combining with the fact that the activation functions are invertible, such constructed \mathbf{g}^* in each trial is also invertible.

We use $N = 20,000$ samples in both domains, denoted as $\{\mathbf{x}_n\}_{n=1}^N, \{\mathbf{y}_n\}_{n=1}^N$, to be the training samples. In addition, we have 1,000 testing samples. The data generation process is as follows:

$$\begin{aligned} \boldsymbol{\mu}_q &\sim \text{Unif}([-1, 1]^2), \quad \forall q \in [Q], \\ \mathbf{x}_{(q-1)N_q+n} &\sim \mathbb{P}_{\mathbf{x}}^{(q)}, \quad \forall n \in [N_q] \quad \forall q \in [Q], \\ \mathbf{y}_{(q-1)N_q+n} &= \mathbf{g}^*(\mathbf{x}_{(q-1)N_q+n}), \end{aligned}$$

where $N_q = \lfloor 20000/Q \rfloor$, indicating that the mixture components have equal probability. In our experiments we use $(\mathbf{x}_n, \mathbf{y}_n)$'s association with one of the Q mixture components as our auxiliary variable. Therefore, we have $I = Q$. In addition, u is uniformly distributed, i.e., $\Pr(u = u_q) = 1/Q, \forall q \in [Q]$, and

$$\mathbb{P}_{\mathbf{x}|u=u_q} = \mathbb{P}_{\mathbf{x}}^{(q)}.$$

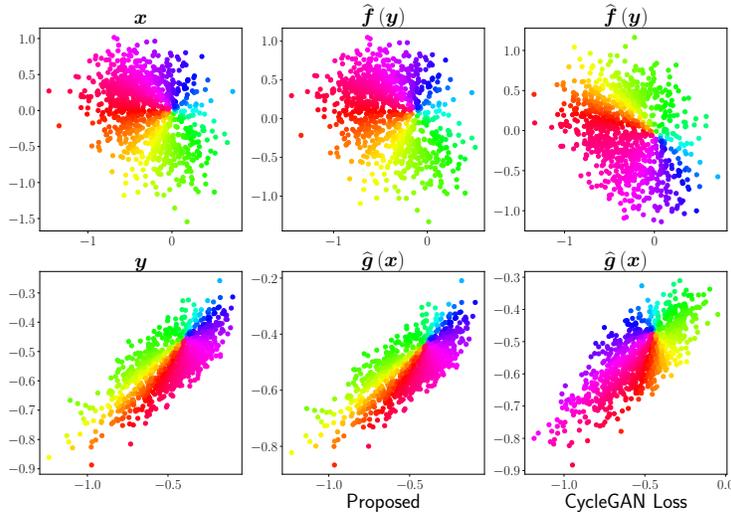


Figure 13: Scatter plots of the source and translated samples. The proposed method uses $I = Q = 3$.

Evaluation Metric. In the synthetic data, we have access to the ground-truth pairs $(\mathbf{x}_n, \mathbf{y}_n)$. Hence, we measure the translation error (TE) using

$$\text{TE} = \sum_{n=1}^N 1/2N (\|\hat{\mathbf{g}}(\mathbf{x}_n) - \mathbf{y}_n\|_2^2 + \|\hat{\mathbf{f}}(\mathbf{y}_n) - \mathbf{x}_n\|_2^2).$$

Implementation Details. To represent $\mathbf{g} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ and $\mathbf{f} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, we use three-layer MLPs, where 256 hidden units are used in each of the 2 hidden layers. We also use leaky ReLU activations with a slope of 0.2. The discriminator is a five-layer MLP with 128 hidden units in each of the hidden layers. Each layer, except for the last, is followed by layer normalization (Ba et al., 2016) and leaky ReLU activations (Maas et al., 2013) with a slope of 0.2. We use the same architecture for all I discriminators in DIMENSION. In the synthetic-data experiments, we implement the distribution matching module using the least-square GAN loss (Mao et al., 2017).

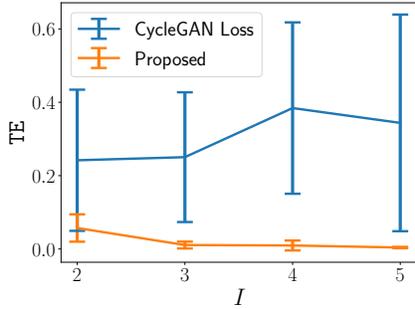
Baseline. In the synthetic experiments, our purpose is to show the lack of translation identifiability of naive distribution matching. Hence, we use the CycleGAN loss in (3) as a benchmark.

Hyperparameter Settings. We use the Adam optimizer with an initial learning rate of 0.0001 with hyperparameters $\beta_1 = 0.5$ and $\beta_2 = 0.999$ (Kingma & Ba, 2015). Note that β_1 and β_2 are hyperparameters of Adam that control the exponential decay rates of first and second order moments, respectively. We use a batch size of 1000 and train the models for 2000 iterations, where one iteration refers to one step of gradient descent of the translation and discriminator neural networks. We use $\lambda = 10$ for (7).

Results. Fig. 13 shows the scatter plots of the original and translated samples for the 1000 testing samples. Here, we set $I = Q = 3$. The original data $\{\mathbf{x}_n\}_{n=1}^N$ and $\{\mathbf{y}_n\}_{n=1}^N$ are plotted on the left-most column. The result of translation using DIMENSION and CycleGAN Loss are presented in the middle and right columns, respectively. In order to qualitatively evaluate the translation performance, we use the same color to plot the paired data points $(\mathbf{x}_n, \mathbf{y}_n)$ and their translations $(\hat{\mathbf{f}}(\mathbf{y}_n), \hat{\mathbf{g}}(\mathbf{x}_n))$. The color is determined by the angle of \mathbf{x}_n in polar coordinates.

As one can see, the supports of \mathbf{x} and $\hat{\mathbf{f}}(\mathbf{y})$ (as well as those of \mathbf{y} and $\hat{\mathbf{g}}(\mathbf{x})$) are well matched by both methods. This implies that both methods can match the distributions fairly well. However, CycleGAN Loss misaligns the samples (by observing the color). The results given by DIMENSION does not have this misalignment issue.

Fig. 14 shows the average TE (over 10 random trials) and the standard deviation attained by DIMENSION and the baseline under different Q 's. Here, we also set $I = Q$ as before. One can see that the average TE decreases with the increase in I . Notably, the *variance* of TE also becomes

Figure 14: TE under various I 's.

much smaller when I grows from 2 to 5—this shows more stable translation performance when I increases. The result is consistent with our theorems, which shows that having a larger I has a better chance to avoid MPA.

F REAL-DATA EXPERIMENT SETTING AND ADDITIONAL RESULTS

In this section we provide details of the real-data experiment settings.

F.1 OBTAINING $\{u_1, \dots, u_I\}$

MrM and ErS Datasets. For these two datasets, we use the available category labels as the alphabet of u . Specifically, for MrM dataset, we use $u \in \{1, \dots, 10\}$, i.e., the labels of the identity of digits. For ErS datasets, we use $u \in \{\text{shoes, sandals, slippers, boots}\}$, which indicates the types of the shoes/edges.

CB Dataset. In this dataset, we designate the alphabet of u to be $u_1 = \text{"black hair"}$, $u_2 = \text{"non-black hair"}$, $u_3 = \text{"male"}$, $u_4 = \text{"female"}$. This information is not fully available in the original CB dataset (to be specific, Bitmoji (Mozafari, 2020) has the gender attributes available but the hair color is not available). We use the foundation model, namely, CLIP (Radford et al., 2021), to acquire the hair color information of each Bitmoji face. Specifically, we use the text prompts “a cartoon of a person whose hair color is mostly black” and “a cartoon of a person whose hair color is not black”. The presence of black hair for each image is decided based on cosine distance of the image embedding with the text embeddings of the two prompts.

F.2 NEURAL NETWORK DETAILS

We use the nomenclature in Table 2 to describe the neural network architecture. For example,

$$\text{Conv}(\text{C-}N_{\text{in}} - N_{\text{out}}, \text{K-}N_k, \text{S-}N_s, \text{ZP-}N_p), \text{LN}, \text{LeakyReLU}$$

refers to a convolutional layer with N_{in} input channels and N_{out} output channels; $\text{K-}N_k$ means that the size of kernel is N_k ; $\text{S-}N_s$ means that the stride is N_s ; and $\text{ZP-}N_p$ means that the zero padding has a size of N_p . The convolutional layer is followed by layer normalization (LN) and then LeakyReLU activations.

The translation neural networks, g and f , for images of size 256×256 follow the architecture outlined in Table 3. For images of size 128×128 , a modified architecture is used, where one down-sampling layer (see Layer #6) and one up-sampling layer (Layer #11) in Table 3 are not included. For images of size 32×32 , three down-sampling layers (indices from #4 to #6) and three up-sampling layers (indices from #11 to #13) are not included.

ResBlock refers to block of convolutional layers with shortcut connection and optional downsampling. Specifically, $\text{ResBlock}(\text{C-}M\text{-}N, \text{Operation})$ is composed of two smaller blocks, namely, $\text{Process}(\text{C-}M\text{-}N, \text{Operation})$ and $\text{Shortcut}(\text{C-}M\text{-}N, \text{Operation})$. The $\text{Process}(\text{C-}M\text{-}N, \text{Operation})$ block has the following layers:

Table 2: Nomenclature for neural network components

Abbreviation	Definition
Conv	Convolutional Layer
IN	Instance normalization
ReLU	ReLU activation
LeakyReLU	Leaky-ReLU activation with 0.2 slope
Tanh	tanh activation function
UpSample	Upsample using nearest neighbor with scale factor of 2
DownSample	Downsample using average pooling with a scale factor of 2
K- N	Kernel (filter) of size N
S- N	Stride of size N
ZP- N	Zero Padding of size N
C- M - N	M input and N output channels

Table 3: Translation neural network architecture for f and g .

Layer Number	Layer Details
1	Conv-(C-3-64, K-1, S-1, ZP-0)
2	ResBlock-(C-64-128, DownSample)
3	ResBlock-(C-128-256, DownSample)
4	ResBlock-(C-256-512, DownSample)
5	ResBlock-(C-512-512, DownSample)
6	ResBlock-(C-512-512, DownSample)
7	ResBlock-(C-512-512, -)
8	ResBlock-(C-512-512, -)
9	ResBlock-(C-512-512, -)
10	ResBlock-(C-512-512, -)
11	ResBlock-(C-512-512, UpSample)
12	ResBlock-(C-512-512, UpSample)
13	ResBlock-(C-512-256, UpSample)
14	ResBlock-(C-256-128, UpSample)
15	ResBlock-(C-128-64, UpSample)
16	Conv-(C-64-3, K-1, S-1, ZP-0)

1. IN, LeakyReLU, Conv-(C- M - M , K-3, S-1, ZP-1)
2. *Operation*
3. IN, LeakyReLU, Conv-(C- M - N , K-3, S-1, ZP-1)

The Shortcut-(C- M - N , *Operation*) block consists of the following layers:

1. Conv-(C- M - N , K-1,S-1,ZP-0)
2. *Operation*

Let z denote the input to the ResBlock and w the output of the ResBlock. Then the forward pass of ResBlock is expressed as follows:

$$w = \text{ResBlock}(z) = \text{Process}(z) + \text{Shortcut}(z).$$

We use multi-task discriminators (Liu et al., 2019) with output dimension of I to represent $\mathbf{d}_x^{(i)}, \mathbf{d}_y^{(i)}, \forall i \in [I]$. Specifically, each of the multi-task discriminators \mathbf{d}_x and \mathbf{d}_y has I output dimensions. The i th outputs of \mathbf{d}_x and \mathbf{d}_y correspond to $\mathbf{d}_x^{(i)}$ and $\mathbf{d}_y^{(i)}$, respectively.

Table 4: Discriminator architecture for $\mathbf{d}_x : \mathbb{R}^{256 \times 256 \times 3} \rightarrow \mathbb{R}^I$, $\mathbf{d}_y : \mathbb{R}^{256 \times 256 \times 3} \rightarrow \mathbb{R}^I$.

Layer Number	Layer Details
1	Conv-(C-3-64, K-1, S-1, ZP-0)
2	ResBlock-(C-64-128, DownSample),
3	ResBlock-(C-128-256, DownSample),
4	ResBlock-(C-256-512, DownSample),
5	ResBlock-(C-512-512, DownSample),
6	ResBlock-(C-512-512, DownSample),
7	ResBlock-(C-512-512, DownSample), LeakyReLU
8	Conv-(C-512,512, K-4, S=1, ZP=0), LeakyReLU
9	Reshape-512
10	Linear-(512,I)

F.3 HYPERPARAMETER SETTING

We use the Adam optimizer with an initial learning rate of 0.0001 with hyperparameters $\beta_1 = 0.0$ and $\beta_2 = 0.999$ (Kingma & Ba, 2015). Note that β_1 and β_2 are hyperparameters of Adam that control the exponential decay rates of first and second order moments, respectively. We set our regularization parameter $\lambda = 10$. We use a batch size of 16. We train the networks for 100,000 iterations. Following standard practice, we add squared ℓ_2 -norm regularization on the network parameters and use a *weight decay* of 0.00001. For the translation tasks with 256×256 images (CelebA-HQ to Bitmoji Faces), the runtime using a single Tesla V100 GPU is approximately 55 hours. For the translation tasks with 128×128 images (Edges to Rotated Shoes), the runtime using a single Tesla V100 GPU is approximately 35 hours. In order to stabilize the GAN training dynamics, we add a gradient penalty term. This term penalizes discriminators’ large gradients, which is known to help the convergence of the GAN objective (Mescheder et al., 2018). We modified the regularization to accommodate our diversified DT loss function. The modified regularization term is as follows:

$$\mathcal{R} = \frac{\gamma}{2} \Pr(u = u_i) \left(\mathbb{E}_{\mathbf{x} \sim \mathcal{P}_{\mathbf{x}|u=u_i}} \|\nabla \mathbf{d}_x^{(i)}\|_2^2 + \mathbb{E}_{\mathbf{y} \sim \mathcal{P}_{\mathbf{y}|u=u_i}} \|\nabla \mathbf{d}_y^{(i)}\|_2^2 \right),$$

where $\nabla \mathbf{d}_x^{(i)}$ denotes the gradient of $\mathbf{d}_x^{(i)}$. We set the value of γ to be 1.0. We take exponential moving average (EMA) of the parameters during training as the final estimate of the parameters of the trained neural networks. We use a weighting factor of 0.999. This has been observed to improve the performance of GANs (Karras et al., 2017; Yaz et al., 2018).

F.4 DATASET DETAILS

MNIST to Rotated MNIST (MrM). We use 60,000 training samples of the MNIST digits (LeCun et al., 2010) that have a dimension 28×28 as the \mathcal{X} -domain. For the \dagger -domain, each of the 60,000 digits is rotated by 90 degrees. The orders of samples are shuffled in both domains to “break” the content correspondence. Under this setting, each \mathbf{x} has a ground-truth correspondence \mathbf{y} .

Edges to Rotated Shoes (ErS). Edges2Shoes dataset (Isola et al., 2017) consists of 49,825 training samples. We resize the all images to have 128×128 pixels. The \mathcal{X} -domain corresponds to the *edges of the shoes*, and the \mathcal{Y} -domain corresponds to the *shoes* that are rotated by 90 degrees. Like in the MrM dataset, the ground-truth correspondence is known to us, which can assist evaluation.

CelebA-HQ to Bitmoji Faces (CB) We use 29,900 training samples from CelebA-HQ (Karras et al., 2017) as the \mathbf{x} -domain, and 3,984 training samples from Bitmoji faces (Mozafari, 2020) as the \mathbf{y} -domain. Note that Bitmoji Faces consists of only 4,084 samples in total, of which 100 samples are held out as the test samples. We resize all images in both domains to have 256×256 pixels. Unlike the previous two datasets, the ground-truth correspondence is *not* known to us in this dataset.

Evaluation Details. The LPIPS score is computed using 100 test samples. Pre-trained AlexNet (Krizhevsky et al., 2012; Zhang et al., 2018) is used in order to compute the LPIPS scores.

The FID score is computed using 1000 translated and real samples for each domain. Pre-trained Inception-v3 (Szegedy et al., 2016) is used in order to compute the FID scores.

F.5 BASELINES

We use `CycleGAN+Id` (Zhu et al., 2017)³, `UNIT` (Liu et al., 2017)⁴, `MUNIT` (Huang et al., 2018)⁴, `U-GAT-IT` (Kim et al., 2020)⁵, `StarGAN-v2` (Choi et al., 2020)⁶, `ZeroDIM` (Gabbay et al., 2021)⁷, `OverLORD` (Gabbay & Hoshen, 2021)⁸, `Hneg-SRC` (Jung et al., 2022)⁹, `GP-UNIT` (Yang et al., 2023)¹⁰, and the plain-vanilla `CycleGAN Loss` in (3) as the baselines.

For `StarGAN-v2` and `GP-UNIT`, training is done with their default settings (specifically, the configurations for the ‘AFHQ’ dataset in their papers are used). For `CycleGAN+Id`, `UNIT`, `MUNIT`, `U-GAT-IT`, and `Hneg-SRC`, we train the models for 200,000 iterations. We use a batch size of 8 for these methods except for `U-GAT-IT`, which uses 4 in order to control the computational load and runtime. These parameters are carefully set for the baselines to our best extent. For `OverLORD` (Gabbay & Hoshen, 2021), we use the setting used for male to female translation task on CelebA-HQ dataset in their paper. For `ZeroDIM` (Gabbay et al., 2021), we use the setting used for experiments on FFHQ dataset, which has a similar size as the datasets used in our paper. Note that `ZeroDIM` also uses the same auxiliary variables as those used in the proposed method.

F.6 ADDITIONAL RESULTS

In this subsection, we present additional qualitative and quantitative results.

Fig. 16 shows the result of translating Bitmoji faces (B) to celebrity profile photos (C). As mentioned in the main text, translating from the B domain to the C domain is a hard task as the learned translation function needs to “fill in” a lot of details to make the generated profiles photorealistic. Visually, one can see that the proposed method (with $I = 4$) exhibits much more intuitive content alignment relative to the baselines. In addition, the proposed method using $I = 2$ (only using ‘male’ and ‘female’ as the auxiliary variable alphabet) also provides more satisfactory results relative to the baselines. This echoes our theoretical claim that the chance of attaining translation identifiability grows quickly when I increases. It also shows that diversifying the distributions to be matched, even if just one more distribution pair is included, helps improve the final performance.

Fig. 15 shows the result of translating edges (E) to rotated shoes (rS). Visually, our method significantly outperforms the baselines in terms of content alignment. It is interesting to notice that, although “edges to shoes” (no rotation) is a well studied dataset, our experiments show that a simple rotation makes most of the existing methods struggle to produce reasonable results. However, our method is insensitive to this kind of geometric changes. In the literature, the baselines `U-GAT-IT` (Kim et al., 2020) and `GP-UNIT` (Yang et al., 2023) were shown to be good at handling certain geometric variations. However, one can see that their performance over the ErS dataset is still far from ideal. The result shows the importance of taking translation identifiability into account, especially when drastic geometric changes happen across domains.

Fig. 18 and Fig. 17 show similar results for the translation of CelebA-HQ (C) to Bitmoji Faces (B) and Rotated Shoes (rS) to Edges (E), respectively. Fig. 19 shows the translations between MNIST (M) and rotated MNIST digits (rM).

³<https://github.com/junyanz/pytorch-CycleGAN-and-pix2pix.git>

⁴<https://github.com/NVlabs/MUNIT.git>

⁵<https://github.com/znxlwm/UGATIT-pytorch.git>

⁶<https://github.com/clovaai/stargan-v2.git>

⁷<https://github.com/avivga/zerodim>

⁸<https://github.com/avivga/overlord>

⁹https://github.com/jcy132/Hneg_SRC.git

¹⁰<https://github.com/williamyang1991/GP-UNIT.git>





Figure 16: Translation of Bitmoji to CelebA-HQ.



Figure 17: Translation of rotated Shoes to Edges.

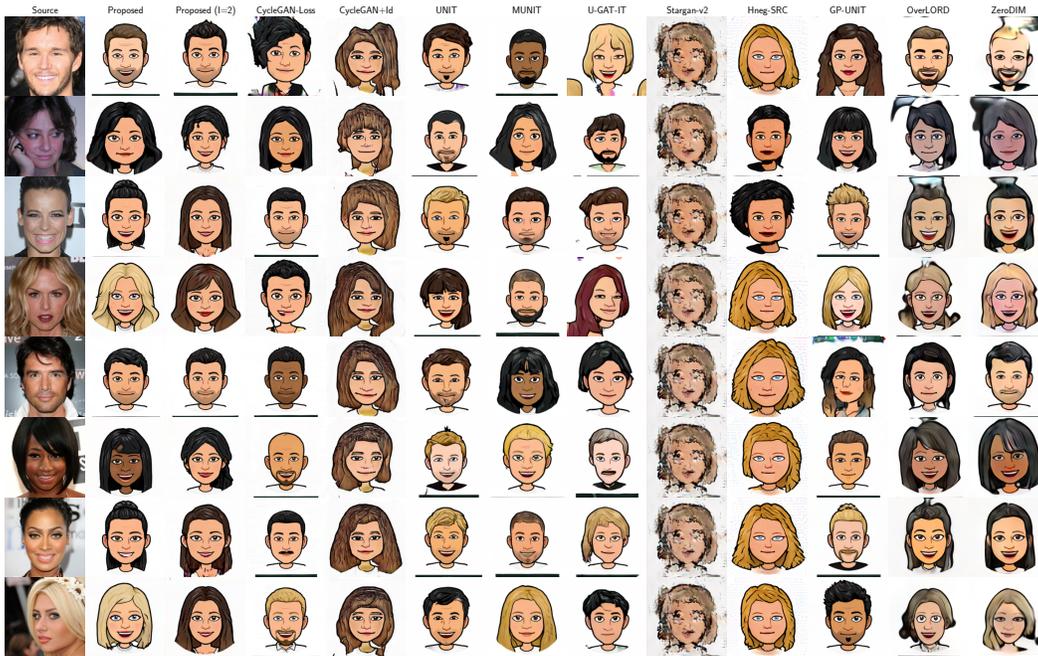


Figure 18: Translation of CelebA-HQ to Bitmoji.

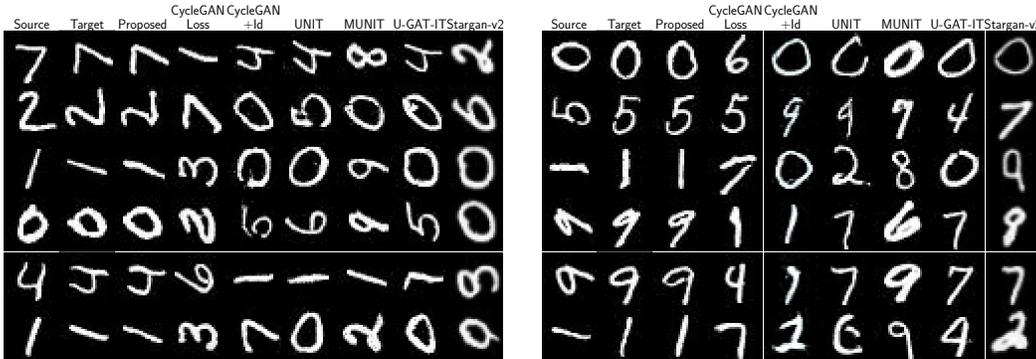


Figure 19: Translation between MNIST digits and rotated MNIST digits.

There are multiple ways to define auxiliary variables for a given UDT task. However, different choices of auxiliary information can result in different level of translation performance. For example, in the MNIST example, one can alternatively use digit shape as the alphabets of the auxiliary variable. Figure 20 shows the result of using different shapes of digits. Here we use the following attributes (with corresponding digits that has those attributes): “line” : [1,2,4,5,7], “circle” : [0,6,8,9], “curve” : [0,2,3,5,6,8,9], “vertical line” : [1, 4, 5], “horizontal line” : [2,5,4,7], “curve without loops” : [2,3,5,6,9], “only vertical line” : [1], “only horizontal line” : [2]. One can see that using digit identity as auxiliary information results in a slightly enhanced performance compared to using the digit shape as auxiliary information.

G IMPROVING EXISTING METHODS USING DIVERSIFIED DISTRIBUTION MATCHING

We hope to emphasize that the diversified distribution matching (DDM) principle can be combined with many other existing UDT approaches to avoid failure cases. In this section, we use our diversified distribution matching module to replace the their original ones in existing paradigms and observe the performance. For the datasets “Edges vs. Rotated Shoes” (ErS) and “CelebA-HQ vs. Bitmoji” (CB), we select the baselines that are able to generate faithful samples in the target domain based on their FID scores (see Table 1. To be specific, for the ErS dataset, we integrate DDM with UNIT (Liu et al., 2017). For the CB dataset, we combine DDM with GP-UNIT (Yang et al., 2023). In both cases, we keep their method-defined regularization terms and other settings unchanged. We refer to the modified methods as UNIT-DDM and GP-UNIT-DDM, respectively.

Our way of combining DDM with these existing approaches is to replace their discriminators. To obtain UNIT-DDM and GP-UNIT-DDM, we modify the discriminator neural networks of UNIT and GP-UNIT into multi-task discriminators. Specifically, for UNIT-DDM, the multi-scale discriminator of UNIT which has one output channel for each scale, is modified to produce I output channels for each scale. Similarly, to obtain GP-UNIT-DDM, the discriminator of GP-UNIT is modified to have I output channels instead of one output channel at the output layer. The i th output channel is interpreted as the i th discriminator associated with u_i .

Fig. 21 shows the qualitative results attained by the original versions of UNIT and GP-UNIT as well as their DDM-modified versions. One can see that there is significant improvement in terms of content alignment, without compromising the visual quality—see the FID and LPIPS scores in Table 5. This attests to the hypothesis that distribution-matching based domain translation frameworks can benefit from the proposed MPA eliminating idea.

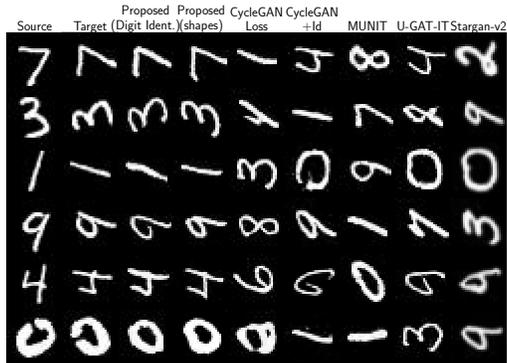


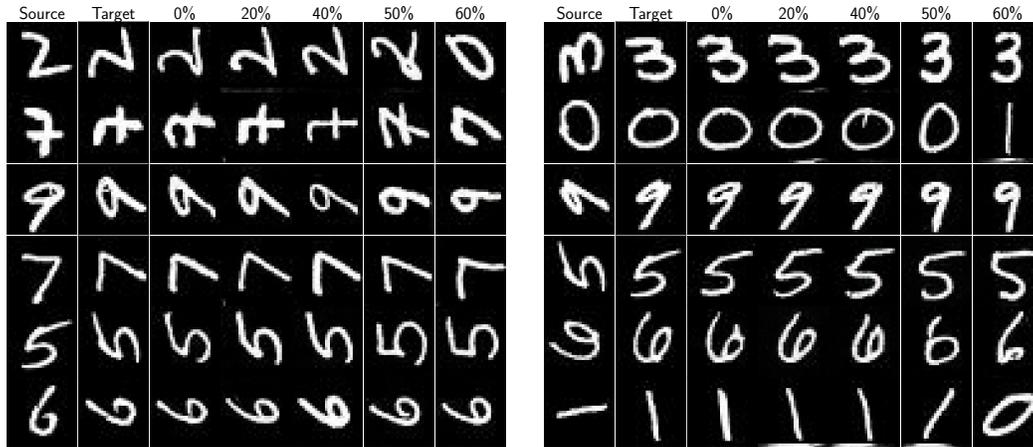
Figure 20: Result of using different auxiliary variable for MNIST digits to rotated MNIST digits task. Using shape attributes incur LPIPS= 0.19 ± 0.08 , compared to LPIPS= 0.11 ± 0.08 for the digit identity as auxiliary variable.

Table 5: The FID and LPIPS scores attained by UNIT and GP-UNIT as well as their DDM-modified versions.

Method	FID (\downarrow)		LPIPS (\downarrow)	
	Edges	Shoes	Edges \rightarrow Rot. Shoes	Rot. Shoes \rightarrow Edges
UNIT	33.95	96.28	0.49 ± 0.035	0.58 ± 0.038
UNIT-DDM	43.95	88.58	0.30 ± 0.075	0.35 ± 0.092
Method	CelebA-HQ	Bitmoji		
GP-UNIT	32.40	30.30		
GP-UNIT-DDM	37.79	30.33		

Table 6: LPIPS score attained by DIMENSION using random u_i assignments.

random u_i proportion	MNIST \rightarrow Rot. MNIST	Rot. MNIST \rightarrow MNIST
0%	0.11 ± 0.082	0.09 ± 0.047
20%	0.09 ± 0.050	0.08 ± 0.040
40%	0.10 ± 0.049	0.13 ± 0.064
50%	0.19 ± 0.080	0.19 ± 0.083
60%	0.25 ± 0.124	0.21 ± 0.086

Figure 22: Result of DIMENSION under random u_i assignments to various fractions of training data.

H ROBUSTNESS TO NOISY AUXILIARY VARIABLES.

It is of interest to know whether using noisy or wrong auxiliary variables would heavily affect the performance of DIMENSION. To this end, we assign random u_i 's to a fraction of the training samples in the "MNIST vs. Rotated MNIST" dataset.

Table 6 and Fig. 22 show the LPIPS scores and qualitative results attained by DIMENSION, respectively, under different fractions of random (and highly possibly wrong) auxiliary variables. Notably, there is almost no performance degradation of DIMENSION even when 40% of the assigned u_i 's are random. This shows the method's robustness to wrong/noisy auxiliary variables.