# TraVeLGAN: Image-to-image Translation by Transformation Vector Learning

Matthew Amodio
Yale University
matthew.amodio@yale.edu

Smita Krishnaswamy
Yale University
smita.krishnaswamy@yale.edu

## Abstract

*Interest in image-to-image translation has grown substantially in recent years with the success of unsupervised models based on the cycle-consistency assumption. The achievements of these models have been limited to a particular subset of domains where this assumption yields good results, namely homogeneous domains that are characterized by style or texture differences. We tackle the challenging problem of image-to-image translation where the domains are defined by high-level shapes and contexts, as well as including significant clutter and heterogeneity. For this purpose, we introduce a novel GAN based on preserving intra-domain vector transformations in a latent space learned by a siamese network. The traditional GAN system introduced a discriminator network to guide the generator into generating images in the target domain. To this two-network system we add a third: a siamese network that guides the generator so that each original image shares semantics with its generated version. With this new three-network system, we no longer need to constrain the generators with the ubiquitous cycle-consistency restraint. As a result, the generators can learn mappings between more complex domains that differ from each other by large differences - not just style or texture).*

## 1. Introduction

Learning to translate an image from one domain to another has been a much studied task in recent years [36, 17, 15, 38, 13]. The task is intuitively defined when we have paired examples of an image in each domain, but unfortunately these are not available in many interesting cases. Enthusiasm has grown as the field has moved towards unsupervised methods that match the distributions of the two domains with generative adversarial networks (GANs) [18, 11, 32, 35, 26]. However, there are infinitely many mappings between the two domains [24], and there is no guarantee that an individual image in one domain will share any characteristics with its representation in the other domain after mapping.
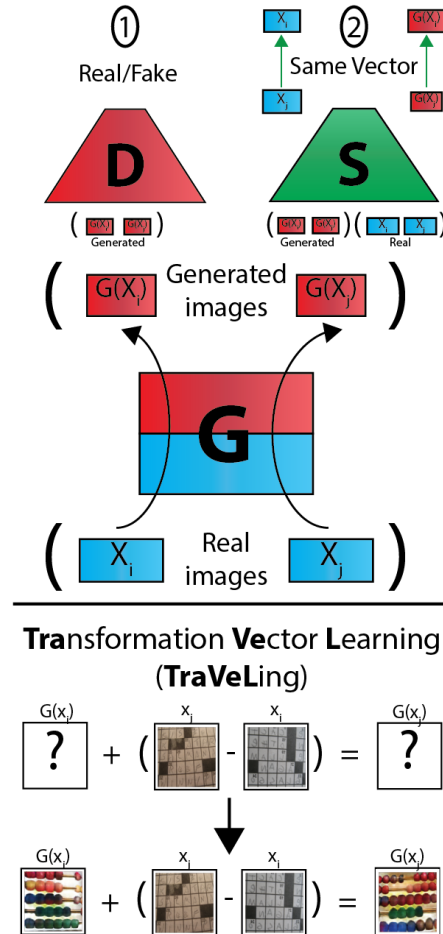


Figure 1: The TraVeLGAN architecture, which adds a siamese network $S$ to the traditional generator $G$ and discriminator $D$ and trains to preserve vector arithmetic between points in the latent space of $S$.

Other methods have addressed this non-identifiability problem by regularizing the family of generators in various ways, including employing cross-domain weight-coupling in some layers [26] and decoding from a shared embedding space [25]. By far the most common regularization, first introduced by the CycleGAN and the DiscoGAN, has been forcing the generators to be each other's inverse, known

Figure 2: Examples of TraVeLGAN generated output on Imagenet domains that are too different and diverse for cycle-consistent GANs to map between. The TraVeLGAN successfully generates images that are both fully realistic in the output domain (shape of object, color, background) and have preserved semantics learned by the siamese network.

as the cycle-consistency property [16, 39, 20, 31, 27, 2, 9, 4, 37]. Recent findings have shown that being able to invert a mapping at the entire dataset level does not necessarily lead to the generation of related real-generated image pairs [23, 3, 11].

Not only do these dataset-level regularizations on the generator not provide individual image-level matching, but also by restricting the generator, they prevent us from learning mappings that may be necessary for some domains. Previous work continues to pile up regularization after regularization, adding restrictions on top of the generators needing to be inverses of each other. These in-

clude forcing the generator to be close to the identity function [39], matching population statistics of discriminator activations [20], weight sharing [26], penalizing distances in the latent space [31], perceptual loss on a previously trained model [25], or more commonly, multiple of these.

Instead of searching for yet another regularization on the generator itself, we introduce an entirely novel approach to the task of unsupervised domain mapping: the **Tra**nsformation **Ve**ctor **L**earning GAN (TraVeLGAN).

The TraVeLGAN uses a third network, a siamese network, in addition to the generator and discriminator to produce a latent space of the data to capture high-level seman-

tics characterizing the domains. This space guides the generator during training, by forcing the generator to preserve vector arithmetic between points in this space. The vector that transforms one image to another in the original domain must be the same vector that transforms the generated version of that image into the generated version of the other image. Inspired by word2vec embeddings [14] in the natural language space, if we need to transform one original image into another original image by moving a foreground object from the top-left corner to the bottom-right corner, then the generator must generate two points in the target domain separated by the same transformation vector.

In word2vec, semantic vector transformations are a *property* of learning a latent space from known word contexts. In TraVeLGAN, we *train* to produce these vectors while learning the space.

Domain mapping consists of two aspects: (a) transfer the given image to the other domain and (b) make the translated image similar to the original image in some way. Previous work has achieved (a) with a separate adversarial discriminator network, but attempted (b) by just restricting the class of generator functions. We propose the natural extension to instead achieve (b) with a separate network, too.

The TraVeLGAN differs from previous work in several substantial ways.

1. It completely eliminates the need for training on cycle-consistency or coupling generator weights or otherwise restricting the generator architecture in any way.

2. It introduces a separate network whose *output* space is used to score similarity between original and generated images. Other work has used a shared latent *embedding* space, but differs in two essential ways: (a) their representations are forced to overlap (instead of preserving vector arithmetic) and (b) the decoder must be able to decode out of the embedding space in an autoencoder fashion [25, 31] ([25] shows this is in fact equivalent to the cycle consistency constraint).

3. It is entirely parameterized by neural networks: nowhere are Euclidean distances between images assumed to be meaningful by using mean-squared error.

4. It adds interpetability to the unsupervised domain transfer task through its latent space, which explains what aspects of any particular image were used to generate its paired image.

As a consequence of these differences, the TraVeLGAN is better able to handle mappings between complex, heterogeneous domains that require significant and diverse shape changing.

By avoiding direct regularization of the generators, the TraVeLGAN also avoids problems that these regularizations cause. For example, cycle-consistency can unnecessarily prefer an easily invertible function to a possibly more coherent one that is slightly harder to invert (or preventing us from mapping to a domain if the inverse is hard to learn). Not only must each generator learn invertible mappings, but it further requires that the two invertible mappings be each other's inverses. Furthermore, cycle-consistency is enforced with a pixel-wise MSE between the original and reconstructed image: other work has identified the problems caused by using pixelwise MSE, such as the tendency to bias towards the mean images [7].

Our approach bears a resemblance to that of the DistanceGAN [6], which preserves pairwise distances between images after mapping. However, they calculate distance directly on the pixel space, while also not preserving any notion of directionality in the space between images. In this paper, we demonstrate the importance of not performing this arithmetic in the pixel space.

Many of these previous attempts have been developed specifically for the task of style transfer, explicitly assuming the domains are characterized by low-level pixel differences (color, resolution, lines) as opposed to high-level semantic differences (shapes and types of specific objects, composition) [7, 37, 13]. We demonstrate that these models do not perform as well at the latter case, while the TraVeLGAN does.

## 2. Model

We denote two data domains $X$ and $Y$, consisting of finite (unpaired) training points $\{x_i\}_{i=1}^{N_x} \in X$ and $\{y_i\}_{i=1}^{N_y} \in Y$, respectively. We seek to learn two mappings, $G_{XY} : X \to Y$ and $G_{YX} : Y \to X$, that map between the domains. Moreover, we want the generators to do more than just mimic the domains at an aggregate level. We want there to be a meaningful and identifiable relationship between the two representations of each point. We claim that this task of unsupervised domain mapping consists of two components: **domain membership** and **individuality**. Without loss of generality, we define these terms with respect to $G_{XY}$ here, with $G_{YX}$ being the same everywhere but with opposite domains.

**Domain membership** The generator should output points in the target domain, i.e. $G_{XY}(X) \in Y$. To enforce this, we use the standard GAN framework of having a discriminator $D_Y$ that tries to distinguish the generator's synthetic output from real samples in $Y$. This yields the typical adversarial loss term $L_{adv}$:

$$L_{adv} = E_X \left[ D_Y(G_{XY}(X)) \right]$$

**Individuality** In addition, our task has a further requirement than just two different points in $X$ each looking like they belong to $Y$. Given $x_i, x_j \in X, i \neq j$, we want there to be some relationship between $x_i$ and $G_{XY}(x_i)$ that justifies why $G_{XY}(x_i)$ is the representation in domain $Y$ for
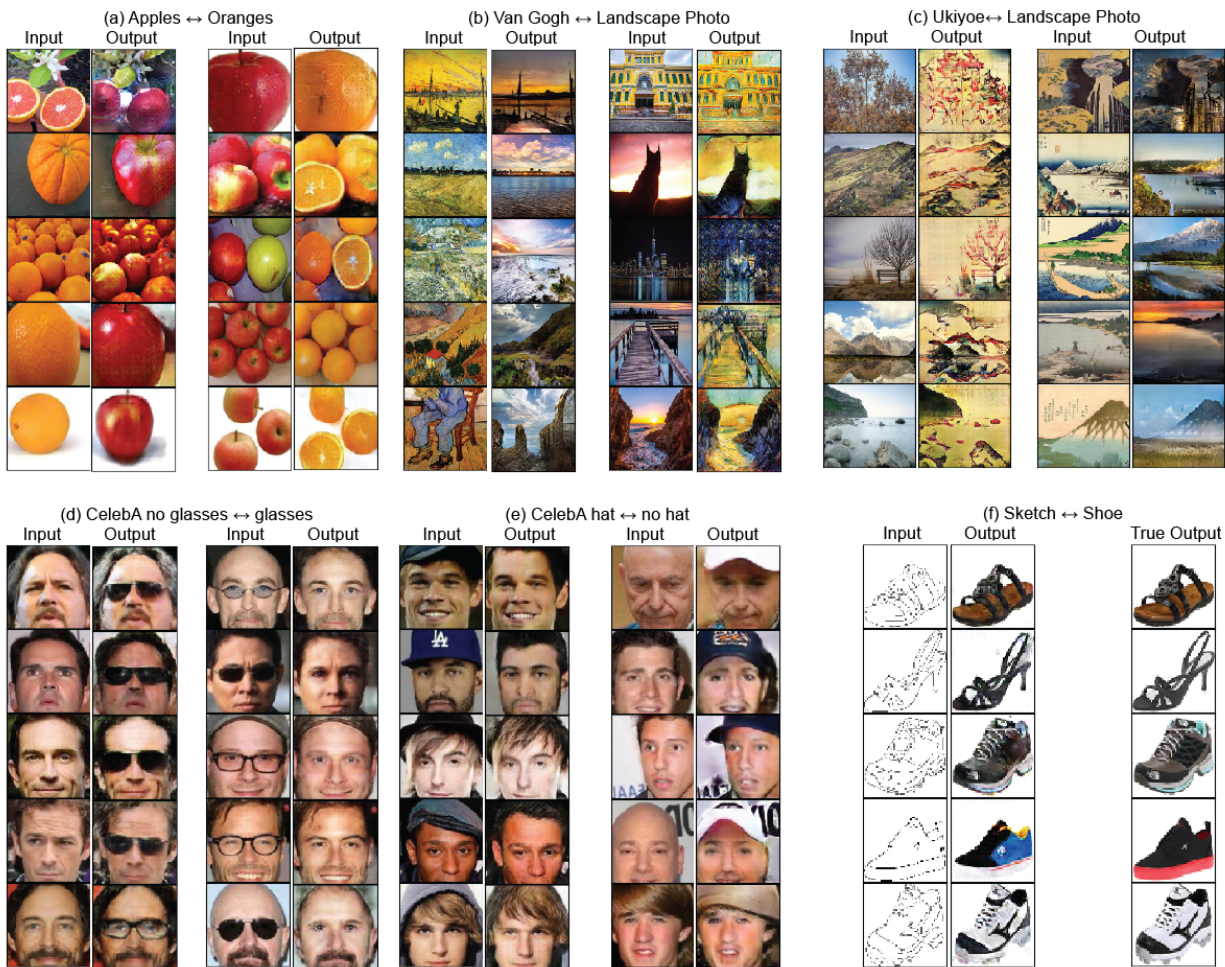
Figure 3: Examples of TraVeLGAN generated output on traditional datasets for unsupervised domain transfer with cycle-consistent GANs. Little change to the original image is necessary in these problems, and TraVeLGAN generates the expected, minimally changed image in the other domain.

$x_i$ and not for $x_j$. Without this requirement, the generator could satisfy its objective by ignoring anything substantive about its input and producing arbitrary members of the other domain.

While other methods try to address this by regularizing $G_{XY}$ (by forcing it to be close to the identity or to be inverted by $G_{YX}$), this limits the ability to map between domains that are too different. So instead of enforcing similarity between the point $x_i$ and the point $G_{XY}(x_i)$ directly in this way, we do so implicitly by matching the relationship between the $x_i$'s and the relationship between the corresponding $G_{XY}(x_i)$'s.

We introduce the notion of a *transformation vector* between two points. In previous natural language processing applications [14], there is a space where the vector that would transform the word *man* to the word *woman* is similar to the vector that would transform *king* to *queen*. In our applications, rather than changing the gender of the word, the

transformation vector could change the background color, size, or shape of an image. The crucial idea, though, is that whatever transformation is necessary to turn one original image into another original image, an analogous transformation must separate the two generated versions of these images.

Formally, given $x_i, x_j \in X$, define the transformation vector between them $\nu(x_i, x_j) = x_j - x_i$. The generator must learn a mapping such that $\nu(x_i, x_j) = \nu(G_{XY}(x_i), G_{XY}(x_j))$. This is a more powerful property than even preserving distances between points, as it requires the space to be organized such that the directions of the vectors as well as the magnitudes be preserved. This property requires that the vector that takes $x_i$ to $x_j$, be the same vector that takes $G_{XY}(x_i)$ to $G_{XY}(x_j)$.

As stated so far, this framework would only be able to define simple transformations, as it is looking directly at the input space. By analogy, the word-gender-changing vec-

tor transformation does not hold over the original one-hot encodings of the words, but instead holds in some reduced semantic latent space. So we instead redefine the transformation vector to be $\nu(x_i, x_j) = S(x_j) - S(x_i)$, where $S$ is a function that gives a representation of each point in some latent space. Given an $S$ that learns high-level semantic representations of each image, we can use our notion of preserving the transformation vectors to guide generation. We propose to learn such a space with an analogue to the adversarial discriminator $D$ from the traditional GAN framework: a cooperative siamese network $S$.

The goal of $S$ is to map images to some space where the relationship between original images is the same as the relationship between their generated versions in the target domain:

$$L_{TraVeL} = \Sigma\Sigma_{i \neq j} Dist(\nu_{ij}, \nu'_{ij})$$
$$\nu_{ij} = S(x_i) - S(x_j)$$
$$\nu'_{ij} = S(G_{XY}(x_i)) - S(G_{XY}(x_j))$$

where $Dist$ is a distance metric, such as cosine similarity. Note this term involves the parameters of $G$, but $G$ needs this space to learn its generative function in the first place. Thus, these two networks depend on each other to achieve their goals. However, unlike in the case of $G$ and $D$, the goals of $G$ and $S$ are not opposed, but cooperative. They both want $L_{TraVeL}$ to be minimized, but $G$ will not learn a trivial function to satisfy this goal, because it also is trying to fool the discriminator. $S$ could still learn a trivial function (such as always outputting zero), so to avoid this we add one further requirement and make its objective multi-task. It must satisfy the standard siamese margin-based contrastive objective [28, 29] $L_{S_c}$, that every point is at least $\delta$ away from every other point in the latent space:

$$L_{S_c} = \Sigma\Sigma_{i \neq j} max(0, (\delta - ||\nu_{ij}||_2))$$

This term incentivizes $S$ to learn a latent space that identifies some differences between images, while $L_{TraVeL}$ incentivizes $S$ to organize it. Thus, the final objective terms of $S$ and $G$ are:

$$L_S = L_{S_c} + L_{TraVeL}$$
$$L_G = L_{adv} + L_{TraVeL}$$

$G$ and $S$ are cooperative in the sense that each is trying to minimize $L_{TraVeL}$, but each has an additional goal specific to its task as well. We jointly train these networks such that together $G$ learns to generate images that $S$ can look at and map to some space where the relationships between original and generated images are preserved.

## 3. Experiments

Our experiments are designed around intentionally difficult image-to-image translation tasks. These translations
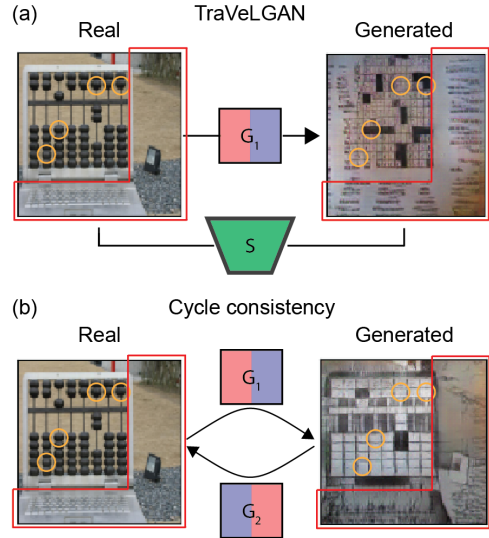


Figure 4: It is hard to learn mappings between domains that are each other's inverse when the domains are asymmetric (e.g. crossword configurations are more complex than abacus configurations). (a) $G_1$ can change the background (red selection) or black beads (orange circles) in hard-to-invert ways. (b) The cycle-consistency assumption forced every black bead to a white crossword square and every blank space to a black crossword square, even though the result is not a realistic crossword pattern. The background is also not fully changed because it could not learn that more complicated inverse function.

are much harder than style or texture transfer problems, where the domain transformation can be reduced to repeating a common patch transformation on every local patch of every image without higher-level semantic information (e.g. turning a picture into a cartoon) [31, 19]. Instead, we choose domains where the differences are higher-level and semantic. For example, when mapping from horses to birds, any given picture of a horse might solely consist of style, texture, and patches that appear in other pictures of real birds (like blue sky, green grass, sharp black outlines, and a brown exterior). Only the higher-level shape and context of the image eventually reveal which domain it belongs to. Additionally, because we use datasets that are designed for classification tasks, the domains contain significant heterogeneity that makes finding commonality within a domain very difficult.

We compare the TraVeLGAN to several previous methods that first regularize the generators by enforcing cycle-consistency and then augment this with further regularizations [39, 20, 3, 31, 27, 2, 9, 4]. Namely, we compare to a GAN with just the cycle-consistency loss (cycle GAN) [39], with cycle-consistency loss plus the identity regularization (cycle+identity GAN) [39], with cycle-consistency loss plus a correspondence loss (cycle+corr GAN) [3], with cycle-consistency loss plus a feature matching regularization (cycle+featmatch GAN) [20], and with

cycle-consistency loss plus a shared latent space regularization (cycle+latent GAN) [25]. The TraVeLGAN utilizes a U-net architecture with skip connections [30] for the generator. The discriminator network is a standard stride-2 convolutional classifier network that doubles the number of filters at each layer until the layer is $4x4$ and outputs a single sigmoidal probability. The siamese network is identical except rather than outputting one node like the discriminator it outputs the number of nodes that is the size of the latent space, without any nonlinear activation. For the cycle-consistent GANs we compare to, we optimized the hyperparameters to get the best achievement we could, since our focus is on testing our different loss formulation. This involved trying both Resnet and U-Net architectures for the models from [39]: the U-Net performed much better than the Resnet at these tasks, so we use that here. We also had to choose a value of the cycle-consistent coefficient that largely de-emphasized it in order to get them to change the input image at all (0.1). Even so, we were not able to achieve nearly as convincing results with any of the baseline models as with the TraVeLGAN.

## 3.1. Similar domains

The datasets we first consider are traditional cases for unsupervised domain mapping with cycle-consistent networks, where little change is necessary. These are:

**Apples to oranges**  The photos of apples and oranges from [39] (Figure 3a). The TraVeLGAN successfully changes not only the color of the fruit, but also the shape and texture. The stem is removed from apples, for example, and the insides of oranges aren't just colored red but fully made into apples. In the last row, the TraVeLGAN changes the shape of the orange to become an apple and correspondingly moves its shadow down in the frame to correspond.

**Van Gogh to landscape photo**  The portraits by Van Gogh and photos of landscapes, also from [39] (Figure 3b). Here the prototypical Van Gogh brush strokes and colors are successfully applied or removed. Notably, in the last row, the portrait of the man is changed to be a photo of a rocky outcrop with the blue clothes of the man changing to blue sky and the chair becoming rocks, rather than becoming a photo-realistic version of that man, which would not belong in the target domain of landscapes.

**Ukiyoe to landscape photo**  Another dataset from [39], paintings by Ukiyoe and photos of landscapes (Figure 3c). It is interesting to note that in the generated Ukiyoe images, the TraVeLGAN correctly matches reflections of mountains in the water, adding color to the top of the mountain and the corresponding bottom of the reflection.

**CelebA glasses**  The CelebA dataset filtered for men with and without glasses [8] (Figure 3d). As expected, the TraVeLGAN produces on the minimal change necessary to trans-

fer an image to the other domain, i.e. adding or removing glasses while preserving the other aspects of the image. Since the TraVeLGAN learns a semantic, rather than pixel-wise, information preserving penalty, in some cases aspects not related to the domain are also changed (like hair color or background). In each case, the resulting image is still a convincingly realistic image in the target domain with a strong similarity to the original, though.

**CelebA hats**  The CelebA dataset filtered for men with and without hats [8] (Figure 3e). As before, the TraVeL-GAN adds or removes a hat while preserving the other semantics in the image.

**Sketch to shoe**  Images of shoes along with their sketch outlines, from [33] (Figure 3f). Because this dataset is paired (though it is still trained unsupervised as always), we are able to quantify the performance of the TraVeLGAN with a heuristic: the pixel-wise mean-squared error (MSE) between the TraVeLGAN's generated output and the true image in the other domain. This can be seen to be a heuristic in the fourth row of Figure 3c, where the blue and black shoe matches the outline of the sketch perfectly, but is not the red and black color that the actual shoe happened to be. However, even as an approximation it provides information. Table 2 shows the full results, and while the vanilla cycle-consistent network performs the best, the TraVeLGAN is not far off and is better than the others. Given that the TraVeLGAN does not have the strict pixel-wise losses of the other models and that the two domains of this dataset are so similar, it is not surprising that the more flexible TraVeL-GAN only performs similarly to the cycle-consistent frameworks. These scores provide an opportunity to gauge the effect of changing the size of the latent space learned by the siamese network. We see that our empirically chosen default value of 1000 slightly outperforms a smaller and lower value. This parameter controls the expressive capability of the model, and the scores suggest providing it too small of a space can limit the complexity of the learned transformation and too large of a space can inhibit the training. The scores are all very similar, though, suggesting it is fairly robust to this choice.

**Quantitative results**  Since the two domains in these datasets are so similar, it is reasonble to evaluate each model using structural similarity (SSIM) between the real and generated images in each case. These results are presented in Table 1. There we can see that the TraVeLGAN performs comparably to the cycle-consistent models. It is expected that the baselines perform well in these cases, as these are the standard applications they were designed to succeed on in the first place; namely, domains that require little change to the original images. Furthermore, it is expected that the TraVeLGAN changes the images slightly more than the models that enforce pixel-wise cycle-consistency. That the

| SSIM | Apple | Van Gogh | Ukiyoe | Glasses | Hats |
|------|-------|----------|--------|---------|------|
| TraVeLGAN | 0.302 | 0.183 | 0.222 | 0.499 | 0.420 |
| Cycle | **0.424** | 0.216 | 0.252 | 0.463 | **0.437** |
| Cycle+ident | 0.305 | **0.327** | **0.260** | **0.608** | 0.358 |
| Cycle+corr | 0.251 | 0.079 | 0.072 | 0.230 | 0.204 |
| Cycle+featmatch | 0.114 | 0.117 | 0.125 | 0.086 | 0.209 |
| Cycle+latent | 0.245 | 0.260 | 0.144 | 0.442 | 0.382 |

Table 1: Real/generated SSIM on the similar-domains datasets.

| Pixel MSE | Sketches | Shoes |
|-----------|----------|-------|
| TraVeLGAN | 0.060 | 0.267 |
| TraVeLGAN ($D_{latent}$=100) | 0.069 | 0.370 |
| TraVeLGAN ($D_{latent}$=2000) | 0.064 | 0.274 |
| Cycle | **0.047** | **0.148** |
| Cycle+corr | 0.427 | 0.603 |
| Cycle+featmatch | 0.077 | 0.394 |
| Cycle+latent | 0.072 | 0.434 |

Table 2: Per-pixel MSE on the shoes-to-sketch dataset.

TraVeLGAN performs so similarly demonstrates quantitatively that the TraVeLGAN can preserve the main qualities of the image when the domains are similar.

### 3.2. Imagenet: diverse domains

The previous datasets considered domains that were very similar to each other. Next, we map between two domains that are not only very different from each other, but from classification datasets where the object characterizing the domain is sometimes only partially in the frame, has many different possible appearances, or have substantial clutter around it. In this most difficult task, we present arbitrary chooses two classes from the Imagenet [10] dataset. These images are much higher-resolution (all images are rescaled to 128x128), making it easier to learn a transfer that only needs local image patches (like style/texture transfer) than entire-image solutions like TraVeLGAN's high-level semantic mappings.

We chose classes arbitrarily because we seek a framework that is flexible enough to make translations between any domains, even when those classes are very different and arbitrarily picked (as opposed to specific domains contrived to satisfy particular assumptions). The pairs are: 1. abacus and crossword (Figure 2a) 2. volcano and jack-o-lantern (Figure 2b) 3. clock and hourglass (Figure 2c) 4. toucan and rock beauty (Figure 2d).

**Asymmetric domains** Learning to map between the domains of abacus and crossword showcase a standard property of arbitrary domain mapping: the amount and nature of variability in one domain is larger than in the other. In Figure 4, we see that the TraVeLGAN learned a semantic mapping from an abacus to a crossword by turning the beads of an abacus into the white squares in a crossword and turning the string in the abacus to the black squares. However, in an abacus, the beads can be aligned in any shape, while in crosswords only specific grids are feasible.
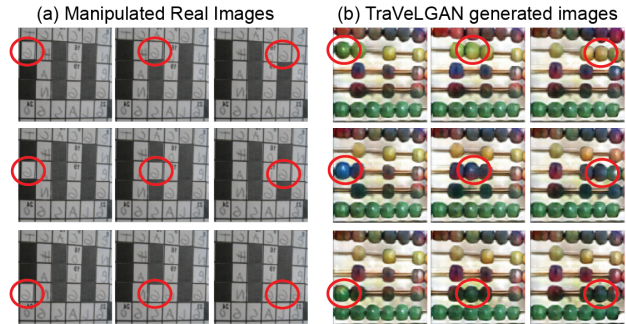


Figure 5: (a) A real crossword image artificially manipulated to move a white square around the frame. (b) The TraVeLGAN, which has not seen any of these images during training, has learned a semantic mapping between the domains that moves an abacus bead appropriately with the crossword square.
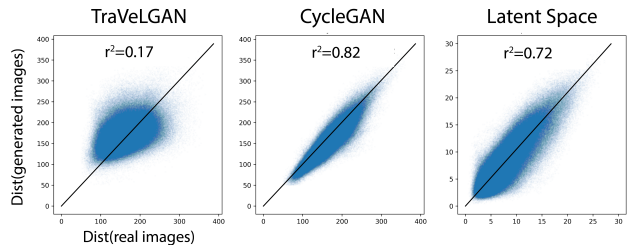


Figure 6: The CycleGAN generates images such that pairwise L2-distances in *pixel space* are strongly preserved. The TraVeLGAN generated images are virtually uncorrelated in pixel space, but the siamese network learns a *latent space* where pairwise distances are preserved.

To turn the abacus in Figure 4 (which has huge blocks of beads that would make for a very difficult crossword indeed!) into a realistic crossword, the TraVeLGAN must make some beads into black squares and others into white squares. The cycle-consistency loss fights this one-to-many mapping because it would be hard for the other generator, which is forced to also be the inverse of this generator, to learn the inverse many-to-one function. So instead, it learns a precise, rigid bead-to-white-square and string-to-black-square mapping at the expense of making a realistic crossword (Figure 4b). Even though the background is an unimportant part of the image semantically, it must recover all of the exact pixel-wise values after cycling. We note that the TraVeLGAN automatically relaxed the one-to-one relationship of beads to crossword squares to create realistic crosswords. On the other hand, any real crossword configuration is a plausible abacus configuration. In the next section, we show that the TraVeLGAN also automatically discovered this mapping can be one-to-one in white-squares-to-beads, and preserves this systematically.

**Manipulated images study** Next we examine the degree to which the TraVeLGAN has learned a meaningful semantic mapping between domains. Since the Imagenet classes

| FID score | (a) | (b) | (c) | (d) |
|---|---|---|---|---|
| TraVeLGAN | **1.026** | **0.032** | **0.698** | **0.206** |
| Cycle | 1.350 | 1.281 | 1.018 | 0.381 |
| Cycle+identity | 1.535 | 0.917 | 1.297 | 1.067 |
| Cycle+corr | 1.519 | 0.527 | 0.727 | 0.638 |
| Cycle+featmatch | 1.357 | 1.331 | 1.084 | 0.869 |
| Cycle+latent | 1.221 | 0.485 | 1.104 | 0.543 |

Table 3: FID scores for each of the models on each of the Imagenet datasets. Column labels correspond to Figure 2.

| Discriminator score | (a) | (b) | (c) | (d) |
|---|---|---|---|---|
| TraVeLGAN | **0.035** | **0.206** | **0.074** | **0.145** |
| Cycle | 0.014 | 0.008 | 0.033 | 0.008 |
| Cycle+identity | 0.011 | 0.044 | 0.040 | 0.064 |
| Cycle+corr | 0.009 | 0.191 | 0.026 | 0.001 |
| Cycle+featmatch | 0.002 | 0.029 | 0.066 | 0.014 |
| Cycle+latent | 0.009 | 0.069 | 0.047 | 0.039 |

Table 4: Discriminator scores for each of the models on each of the Imagenet datasets. Column labels correspond to Figure 2.

are so cluttered and heterogeneous and lack repetition in the form of two very similar images, we create similar images with a manipulation study. We have taken one of the real images in the crossword domain, and using standard photo-editing software, we have created systematically related images. With these systematically related images, we can test to see whether the TraVeLGAN's mapping preserves the semantics in the abacus domain in a systematic way, too.

In Figure 5, we started with a crossword and created a regular three-by-three grid of black squares by editing an image from Figure **??**. Then, systematically, we move a white square around the grid through each of the nine positions. In each case, the TraVeLGAN generates an abacus with a bead moving around the grid appropriately. Remarkably, it even colors the bead to fit with the neighboring beads, which differ throughout the grid. Given that none of the specific nine images in Figure 5 were seen in training, the TraVeLGAN has clearly learned the semantics of the mapping rather than memorizing a specific point.

**Pairwise distance preservation** The DistanceGAN [6] has shown that approximately maintaining pairwise distances between images in ***pixel space*** achieves similar success to the cycle-consistent GANs. In fact, they show that cycle-consistent GANs produce images that preserve the pixel pairwise distance between images with extremely highly correlation. On the toucan to rock beauty dataset, we observe the same phenomenon ($r^2 = 0.82$ in Figure 6). While this produced plausible images in some cases, maintaining pixel-wise distance between images could not generate realistic toucans or rock beauties. The TraVeLGAN pairwise distances are virtually uncorrelated in pixel space ($r^2 = 0.17$). However, we understand the role of the siamese network when we look at the pairwise distances between real images in ***latent space*** and the correspond-

ing pairwise distances between generated images in ***latent space***. There we see a similar correlation ($r^2 = 0.72$). In other words, the TraVeLGAN simultaneously learns a mapping with a neural network to a space where distances can be meaningfully preserved while using that mapping to guide it in generating realistic images.

**Quantitative results** Lastly, we add quantitative evidence to the qualitative evidence already presented that the TraVeLGAN outperforms existing models when the domains are very different. While we used the SSIM and pixel-wise MSE in the previous sections to evaluate success, neither heuristic is appropriate for these datasets. The goal in these mappings is *not* to leave the image unchanged and as similar to the original as possible, it is to fully change the image into the other domain. Thus, we apply use two different metrics to evaluate the models quantitatively on these Imagenet datasets.

In general, quantifying GAN quality is a hard task [5]. Moreover, here we are specifically interested in how well a generated image is paired or corresponding to the original image, point-by-point. To the best of our knowledge, there is no current way to measure this quantitatively for arbitrary domains, so we have pursued the qualitative evaluations in the previous sections. However, in addition to those qualitative evaluation of the correspondence aspect, we at least quantify how well the generated images resemble the target domain, at a population level, with heuristic scores designed to measure this under certain assumptions. The first, the Fréchet Inception Distance (FID score) [12] is an improved version of the Inception Score (whose flaws were well articulated in [5]) which compares the real and generated images in a layer of a pre-trained Inception network (Table 3). The second, the discriminator score, trains a discriminator from scratch, independent of the one used during training, to try to distinguish between real and generated examples (Table 4). The TraVeLGAN achieved better scores than any of the baseline models with both metrics and across all datasets.

## 4. Discussion

In recent years, unsupervised domain mapping has been dominated by approaches building off of the cycle-consistency assumption and framework. We have identified that some cluttered, heterogeneous, asymmetric domains cannot be successfully mapped between by generators trained on this cycle-consistency approach. Further improving the flexibility of domain mapping models may need to proceed without the cycle-consistent assumption, as we have done here.

# References

[1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: a system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016. 11

[2] A. Almahairi, S. Rajeswar, A. Sordoni, P. Bachman, and A. Courville. Augmented cyclegan: Learning many-to-many mappings from unpaired data. *arXiv preprint arXiv:1802.10151*, 2018. 1, 5

[3] M. Amodio and S. Krishnaswamy. Magan: Aligning biological manifolds. *arXiv preprint arXiv:1803.00385*, 2018. 2, 5, 11

[4] A. Anoosheh, E. Agustsson, R. Timofte, and L. Van Gool. Combogan: Unrestrained scalability for image domain translation. *arXiv preprint arXiv:1712.06909*, 2017. 1, 5

[5] S. Barratt and R. Sharma. A note on the inception score. *arXiv preprint arXiv:1801.01973*, 2018. 8

[6] S. Benaim and L. Wolf. One-sided unsupervised domain mapping. In *Advances in neural information processing systems*, pages 752–762, 2017. 3, 8

[7] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 7, 2017. 3

[8] Large-scale celebfaces attributes (celeba) dataset. http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html. Accessed: 2018-10-20. 6

[9] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. *arXiv preprint*, 1711, 2017. 1, 5

[10] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. Ieee, 2009. 7

[11] V. Dumoulin, I. Belghazi, B. Poole, O. Mastropietro, A. Lamb, M. Arjovsky, and A. Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016. 1, 2

[12] Fréchet inception distance (fid score) in pytorch. https://github.com/mseitzer/pytorch-fid. Accessed: 2018-10-20. 8, 11

[13] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2414–2423. IEEE, 2016. 1, 3

[14] Y. Goldberg and O. Levy. word2vec explained: deriving mikolov et al.'s negative-sampling word-embedding method. *arXiv preprint arXiv:1402.3722*, 2014. 3, 4

[15] J. Hoffman, E. Rodner, J. Donahue, T. Darrell, and K. Saenko. Efficient learning of domain-invariant image representations. *arXiv preprint arXiv:1301.3224*, 2013. 1

[16] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. A. Efros, and T. Darrell. Cycada: Cycle-consistent adversarial domain adaptation. *arXiv preprint arXiv:1711.03213*, 2017. 1

[17] X. Huang and S. J. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, pages 1510–1519, 2017. 1

[18] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. *arXiv preprint*, 2017. 1

[19] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711. Springer, 2016. 5

[20] T. Kim, M. Cha, H. Kim, J. K. Lee, and J. Kim. Learning to discover cross-domain relations with generative adversarial networks. *arXiv preprint arXiv:1703.05192*, 2017. 1, 2, 5, 11

[21] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 11

[22] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009. 11

[23] C. Li, H. Liu, C. Chen, Y. Pu, L. Chen, R. Henao, and L. Carin. Alice: Towards understanding adversarial learning for joint distribution matching. In *Advances in Neural Information Processing Systems*, pages 5495–5503, 2017. 2

[24] T. Lindvall. *Lectures on the coupling method*. Courier Corporation, 2002. 1

[25] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *Advances in Neural Information Processing Systems*, pages 700–708, 2017. 1, 2, 3, 6

[26] M.-Y. Liu and O. Tuzel. Coupled generative adversarial networks. In *Advances in neural information processing systems*, pages 469–477, 2016. 1, 2

[27] Y. Lu, Y.-W. Tai, and C.-K. Tang. Conditional cyclegan for attribute guided face image generation. *arXiv preprint arXiv:1705.09966*, 2017. 1, 5

[28] I. Melekhov, J. Kannala, and E. Rahtu. Siamese network features for image matching. In *Pattern Recognition (ICPR), 2016 23rd International Conference on*, pages 378–383. IEEE, 2016. 5

[29] E.-J. Ong, S. Husain, and M. Bober. Siamese network of deep fisher-vector descriptors for image retrieval. *arXiv preprint arXiv:1702.00338*, 2017. 5

[30] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 6, 11

[31] A. Royer, K. Bousmalis, S. Gouws, F. Bertsch, I. Moressi, F. Cole, and K. Murphy. Xgan: Unsupervised image-to-image translation for many-to-many mappings. *arXiv preprint arXiv:1711.05139*, 2017. 1, 2, 3, 5

[32] P. Russo, F. M. Carlucci, T. Tommasi, and B. Caputo. From source to target and back: symmetric bi-directional adaptive gan. *arXiv preprint arXiv:1705.08824*, 2017. 1

[33] igan. https://github.com/junyanz/iGAN/tree/master/train_dcgan. Accessed: 2019-02-01. 6

[34] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015. 11

[35] Y. Taigman, A. Polyak, and L. Wolf. Unsupervised cross-domain image generation. *arXiv preprint arXiv:1611.02200*, 2016. 1

[36] T. Yao, Y. Pan, C.-W. Ngo, H. Li, and T. Mei. Semi-supervised domain adaptation with subspace learning for visual recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2142–2150, 2015. 1

[37] Z. Yi, H. R. Zhang, P. Tan, and M. Gong. Dualgan: Unsupervised dual learning for image-to-image translation. In *ICCV*, pages 2868–2876, 2017. 1, 3

[38] W. Zhang, C. Cao, S. Chen, J. Liu, and X. Tang. Style transfer via image component analysis. *IEEE Transactions on multimedia*, 15(7):1594–1601, 2013. 1

[39] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint*, 2017. 1, 2, 5, 6, 11

**Quantitative results** Quantitative results are summarized by the FID score (Table 3) and the discriminator score (Table 4). We note that these scores were both designed to evaluate models that attempt to generate the full diversity of the Imagenet dataset, while in our case we only map to a single class.

The Fréchet Inception Distance (FID score) [12] calculates the Fréchet distance between Gaussian models of the output of a the pre-trained Inception network [34] on real and generated images, respectively. Lower distances indicate better performance. The results are the mean of the scores from each direction.

The discriminator score is calculated by training a new discriminator, distinct from the one used during training, to distinguish between real and generated images in a domain. A score of zero means the discriminator was certain every generated image was fake, while higher scores indicate the generated images looked more like real images. As in the FID, the results are the mean of the scores from each direction.

**Optimization and training parameters** Optimization was performed with the adam [21] optimizer with a learning rate of 0.0002, $\beta_1 = 0.5$, $\beta_2 = 0.9$. Gradient descent was alternated between generator and discriminator, with the discriminator receiving real and generated images in distinct batches.

**Architecture** The TraVeLGAN architecture is as follows. Let $d$ denote the size of the image. Let $c_n$ be a standard stride-two convolutional layer with $n$ filters, $t_n$ be a stride-two convolutional transpose layer with kernel size four and $n$ filters, and $f_n$ be a fully connected layer outputting $n$ neurons. The discriminator $D$ has layers until the size of the input is four-by-four, increasing the number of filters by a factor of two each time, up to a maximum of eight times the original number (three layers for CIFAR and five layers for Imagenet). This last layer is then flattened and passed through a fully connected layer. The overall architecture is thus $c_n - c_{2n} - c_{4n} - c_{8n} - c_{8n} - f_1$. The siamese network has the same structure as the discriminator except its latent space has size 1000, yielding the architecture $c_n - c_{2n} - c_{4n} - c_{8n} - c_{8n} - f_{1000}$. The generator uses the U-Net architecture [30] that has skip connections that concatenate the input in the symmetric encoder with the decoder, yielding layers of $c_n - c_{2n} - c_{4n} - c_{4n} - c_{4n} - t_{8n} - t_{8n} - t_{8n} - t_{4n} - t_{2n} - t_3$. For the cycle-consistency networks, the architectures of the original implementations were used, with code from [39], [39], [3], [20], for the cycle, cycle+identity, cycle+corr, and cycle+featmatch, respectively. All activations are leaky rectified linear units with leak of 0.2, except for the output layers, which use sigmoid for the discriminator, hyperbolic tangent for the gen-

erator, and linear for the siamese network. Batch normalization is used for every layer except the first layer of the discriminator. All code was implemented in Tensorflow [1] on a single NVIDIA Titan X GPU.

**CIFAR** While the CIFAR images [22] are relatively simple and low-dimensional, it is a deceptively complex task compared to standard domain mapping datasets like CelebA, where they are all centered close-ups of human faces (i.e. their shoulders or hair are in the same pixel locations). The cycle-consistent GANs struggle to identify the characteristic shapes of each domain, instead either only make small changes to the images or focusing on the color tone. The TraVeLGAN, on the other hand, fully transfers images to the target domain. Furthermore, the TraVeLGAN preserves semantics like orientation, background color, body color, and composition in the pair of image (complete comparison results in Figure **??**)

**Interpretability** As the siamese latent space is learned to preserve vector transformations between images, we can look at how that space is organized to tell us what transformation the network learned at a dataset-wide resolution. Figure S1 shows a PCA visualization of the siamese space of the CIFAR dataset with all of the original domain one (bird) and domain two (horse) images. There we can see that $S$ learned a logical space with obvious structure, where mostly grassy images are in the bottom left, mostly sky images in the top right, and so forth. Furthermore, the layout is analogous between the two domains, verifying that the network automatically learned a notion of similarity between the two domains. We also show every generated image across the whole dataset in this space, where we see that the transformation vectors are not just interpretable for some individual images and not others, but are interpretable across the entire distribution of generated images.

**Salience** We next perform a salience analysis of the TraVeL loss by calculating the magnitude of the gradient at each pixel in the generated image with respect to each pixel in the original image (Figure S2). Since the TraVeL loss, which enforces the similarity aspect of the domain mapping problem, is parameterized by another neural network $S$, the original image contributes to the generated image in a complex, high-level way, and as such the gradients are spread richly over the entire foreground of the image. This allows the generator to make realistic abacus beads, which need to be round and shaded, out of square and uniform pixels in the crossword. By contrast, the cycle-consistency loss requires numerical precision in the pixels, and as such the salience map largely looks like a grayscale version of the real image, with rigid lines and large blocks of homogeneous pixels still visible. This is further evidence that the cycle-consistency
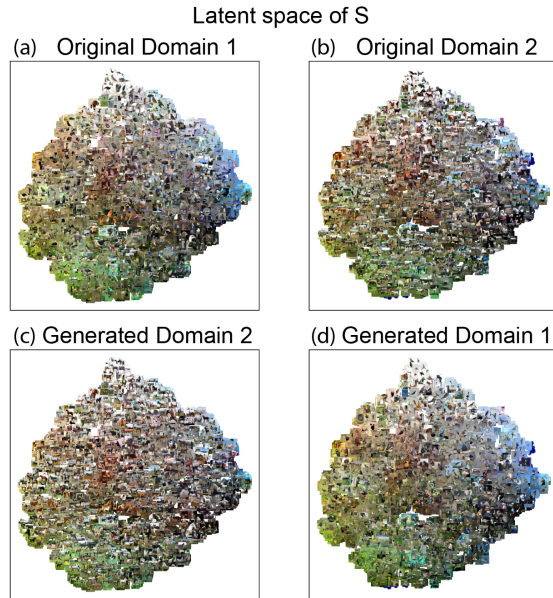
Latent space of S

(a) Original Domain 1

(b) Original Domain 2



(c) Generated Domain 2

(d) Generated Domain 1



Figure S1: Having access to the siamese space output by $S$ provides an interpretability of the TraVeLGAN's domain mapping that other networks lack. PCA visualizations on the CIFAR example indicate $S$ has indeed learned a meaningfully organized space for $G$ to preserve transformation vectors within.

Original Images



TraVeLGAN          Cycle consistency
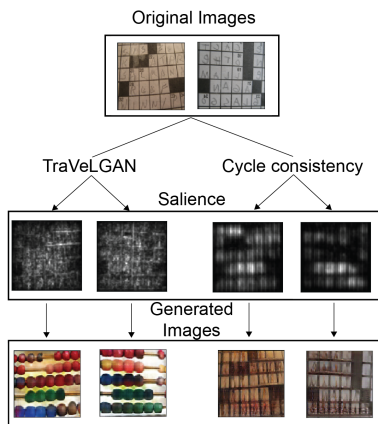
Salience

Generated Images

Figure S2: A salience analysis exploring how the TraVeLGAN's objective loosens the restriction of cycle-consistency and allows it more flexibility in changing the image during domain transfer. The TraVeL loss requires significantly less memorization of the input pixels, and as a result, more complex transformations can be learned.

loss is preventing the generator from making round beads with colors that vary over the numerical RGB values.