# Unpaired Image Translation via Vector Symbolic Architectures

Justin Theiss[1,2], Jay Leverett[1], Daeil Kim[1], and Aayush Prakash[1]

[1] Meta Reality Labs
{theiss,jayleverett,daeilkim,aayushp}@fb.com
[2] University of California, Berkeley, CA 94720, USA

**Abstract.** Image-to-image translation has played an important role in enabling synthetic data for computer vision. However, if the source and target domains have a large semantic mismatch, existing techniques often suffer from source content corruption aka semantic flipping. To address this problem, we propose a new paradigm for image-to-image translation using Vector Symbolic Architectures (VSA), a theoretical framework which defines algebraic operations in a high-dimensional vector (hypervector) space. We introduce VSA-based constraints on adversarial learning for source-to-target translations by learning a hypervector mapping that inverts the translation to ensure consistency with source content. We show both qualitatively and quantitatively that our method improves over other state-of-the-art techniques.

**Keywords:** Image-to-image translation, adversarial learning, vector symbolic architectures, semantic flipping.

## 1 Introduction

Image-to-image translation techniques [8, 12, 14, 23, 25, 30, 33, 39] have been instrumental in improving synthetic data for computer vision. They have been used to bridge the domain gap for synthetic data [12, 30, 39] and for photo-realistic enhancement in virtual reality and gaming applications [33]. Some researchers [7, 16, 31] argue that the domain gap can be further factorized into a content (shift in semantic statistics) and appearance gap. In this work, we are interested in the unpaired image-to-image translation method in the challenging scenario where the content gap between the source and target domains is large.

There are several techniques for unpaired image-to-image translation [12, 30, 39]. Some approaches assume that content and style can be separated in order to translate style without corrupting content (e.g., [12]). Others have used a bijective mapping to reconstruct the source image from the translated image (i.e., a "cyclic loss" [39]). However, these methods do not work well if there is a *large shift in distribution* between source and target domains leading to
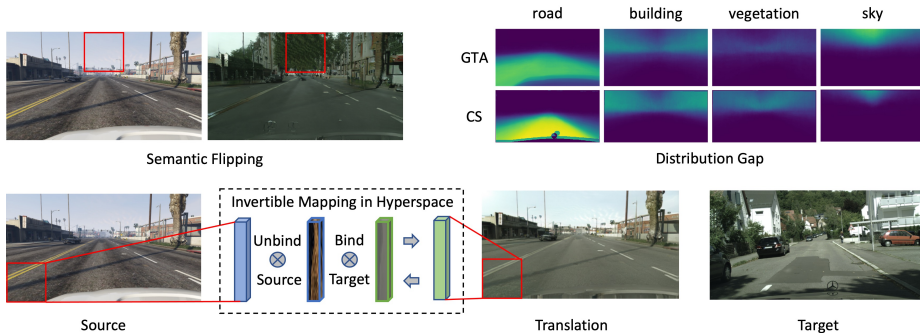
https://github.com/facebookresearch/vsait

Fig. 1: We propose Vector Symbolic Architecture based image-to-image translation technique (VSAIT) which addresses semantic flipping (source content corruption) that happens when the distribution gap (shift in semantic statistics) between source and target domains is large. Our method learns a mapping in high-dimensional vector space (hyperspace), which encourages translated images to be consistent with the source domain. Conceptually, our approach aims to "unbind" source-related information (e.g., texture and color) and "bind" target-related information as well as vice versa to recover source content.

the problem of *semantic flipping*. Semantic flipping is characterized by image artifacts, object or feature hallucinations that are a result of adversarial training for datasets with a large content gap. Specifically, this is observed as a change in content between source and translated images (e.g., sky to trees in Figure 1). Semantic flipping is a critical issue for improving photorealism in computer graphics applications [33] as well as training downstream tasks using translated images, as we want to preserve the source semantic labels of translations.

Relatively few works have directly focused on the semantic flipping problem; however, they can be broadly categorized into three approaches: image-level consistency, domain-invariant encoding, and task-level consistency. Methods using an image-level consistency loss attempt to ensure that pixel-level information is highly correlated between source and translated images [3, 8]. Such methods may fail to account for feature-level differences that do not correlate well with pixel-level differences. Approaches that focus on domain-invariant encoding train the generator to encode domain-invariant content either using a shared latent space [23] or contrastive learning [14]. These methods will fail to reduce semantic flipping if content and style cannot be sufficiently disentangled. Finally, others have used pre-trained task networks to generate pseudo-labels for each domain to ensure a consistent translation with respect to source labels (i.e., semantic masks) [11, 33]. However, these methods fail if the task network cannot generate pseudo-labels for both domains. We instead focus on a method that provides feature-level consistency between source and translated images without explicit assumptions of domain-invariant encoding or quality of pseudo-labels.

This addresses gaps in previous methods which make them vulnerable to semantic flipping.

In the current paper, we propose a novel usage of a theoretical framework known as vector symbolic architectures (VSA [9, 15]; also referred to as hyperdimensional computing) as a new paradigm for unpaired image-to-image translation. Although much of the VSA research has been conducted in theoretical neuroscience (see [21, 20]), it has more recently been applied to computer vision problems using extracted features from deep neural networks [27, 28]. VSA defines a high-dimensional vector (hypervector) space and operators used for symbolic computation, which allows for mathematical formulations of conceptual queries such as, "what color is the car?". In the case of unpaired image translation, such formulations are useful because they enable us to recover attributes from the source image and ensure consistent relationships among features when translating images from one domain to another. VSA is well suited for this approach as it can represent arbitrary symbols and formulations generally without supervision or training [32]. The important difference from previous methods addressing semantic flipping is that VSA ensures that hypervector representations of different semantic content are almost orthogonal to each other (e.g., sky and trees) [15]. The cosine distance between translated and source hypervectors therefore is greatest when semantic flipping occurs.

Using this framework, we propose a method for learning a hypervector mapping between source and target domains that is robust to semantic flipping (Figure 1). By inverting this mapping, we are able to minimize the distance between features extracted from source and translated images without requiring that content and style be fully disentangled. We demonstrate qualitatively and quantitatively that this approach significantly reduces the image artifacts and hallucinations observed for unpaired image translation between domains with a large content gap. We hope that our work provides inspiration for incorporating VSA into new areas of computer vision research. Our contributions include:

- Our method addresses important artifacts and feature hallucinations (semantic flipping) that often occur with other unsupervised image translation methods as a result of content gap between source and target domains.
- To the best of our knowledge, we are the first to show that Vector Symbolic Architectures (VSA) can be used for a challenging task of image translations in the wild. We hope this opens up an exciting area of research around VSA.
- We demonstrate qualitative and quantitative improvement over state-of-the-art methods that directly address semantic flipping across multiple experiments with unmatched semantic statistics.

## 2   Related Work

**Unpaired Image-to-Image Translation** Unpaired image-to-image translation is a widely studied problem with many existing techniques [2, 3, 8, 12, 13, 14, 18, 23, 25, 30, 33, 38, 39]. One popular approach is to impose the
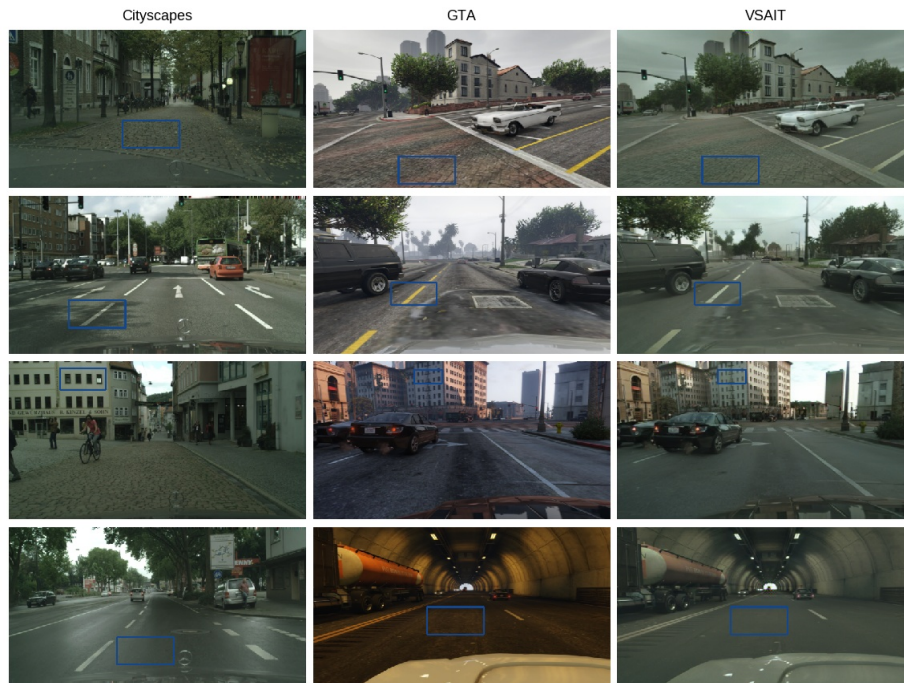
Fig. 2: Example GTA translations from our method (VSAIT) alongside representative examples from Cityscapes. We can see that VSAIT successfully translates GTA images to Cityscapes' style and that, in particular, our method is able to learn specific textures and attributes (e.g. gray cobblestone, white lane lines rather than yellow) that are commonly found in Cityscapes images. Moreover, VSAIT is able to handle unseen scenarios (e.g. tunnel) that do not occur in Cityscapes without semantic flipping.

cycle-consistency constraint [18, 38, 39], which states that a source-to-target generator and a target-to-source generator are bijections. Building on cycle-consistency, UNIT [25] mapped the source and target domains to a shared latent space, while DRIT [23] and other methods [2, 12] factorized the latent space into a shared content space and a domain-specific style space. Distance-GAN [3] showed that source to target translations can be learnt unidirectionally. CUT [30] subsequently proposed a unidirectional method based on contrastive learning. Whereas some of these methods require separate networks per domain or assume that content and style can be disentangled, our method uses a single generator without assumption of disentanglement.

**Semantic Flipping** The body of research on semantic flipping is relatively small. Methods based on cycle-consistency combat semantic flipping through pixel-level reconstruction, but they are unable to prevent semantic flipping in

the face of large distributional shifts. GcGAN [8] enforces geometry-consistent constraints on image translations by ensuring robustness for transformations such as flipping or rotation. This does not necessarily address semantic flipping for small, local features that could still be geometry-consistent. EPE [33] addresses semantic flipping by conditioning on G-buffers (e.g. albedo, glossiness, depth, normals) as well as incorporating semantic segmentation pseudo-labels into the discriminator. Although their method addresses semantic flipping, it does not generalize to many tasks or datasets since it assumes access to the synthetic image rendering process and a downstream task network that can generate pseudo-labels for both source and target domains. Alternatively, SRUNIT [14] proposed to address semantic flipping by using contrastive learning (following CUT [30]) to encode domain-invariant semantic representations that are robust to small feature-space perturbations. Although their approach does not require access to any labels or rendering process like EPE, it does not sufficiently address semantic flipping. Unlike these approaches, we reduce semantic flipping by inverting the translation in hyperspace using VSA in order to ensure recovery of source information.

**Computer Vision Applications of VSA**  Although much of the research associated with VSA has been conducted in theoretical neuroscience, Neubert et al. [28] provided examples of how VSA can be used for computer vision tasks such as object and place recognition, which demonstrated how image features extracted from pre-trained neural networks can be used within the VSA framework. Other works have used VSA for visual question answering [26] and scene transformation [17], albeit with simple shapes or MNIST digits. More recently, Osipov et al. [29] used VSA to learn to unbind category representations among unlabelled data vectors bound to a shared representational space. This is most similar to the challenge that we are addressing in the current paper, as we can consider the source and target features to be bound to a similar content space; however, unlike their study we do not have access to the underlying shared representational space.

## 3   Method

### 3.1   Preliminary: VSA

Vector Symbolic Architectures (VSA) [9, 15] provide a framework for encoding and manipulating information in a hypervector space (hyperspace) with high capacity and robustness to noise [15]. In VSA framework, hypervectors are randomly sampled from a vector space $\mathbb{V} = [-1, 1]^n$ (where $n >> 1000$) in order to represent symbols that can be typically manipulated using two operators – binding/unbinding (element-wise multiplication) and bundling (element-wise addition). The binding operator can be used to assign an attribute to a symbol (e.g., "gray car"), whereas bundling can be used to group symbols together (e.g.,

"gray car and red bike"). These operators define the Multiplication-Addition-Permutation (MAP) architecture. Other operators and architectures (see [35]) are beyond the scope of our work.

The benefit of using VSA framework for image translation is its versatility in representing arbitrary symbols without supervision. This allows us to infer underlying representational spaces (e.g., content distribution) without explicitly defining them. Since the challenge of semantic flipping in unpaired image translation is typically to constrain the generator, we do not need to learn how to generate images from hypervectors but instead learn a mapping that ensures we recover the same source information for a given translated hypervector.

In order to assign random hypervectors to image features, Neubert et al. [28] used locality sensitive hashing (LSH) to project image features extracted using a pre-trained neural network (AlexNet [22]) to a relatively lower-dimensional space (from $13 \times 13 \times 384$ to 8192). In this case, the entire image was encoded into a single hypervector; however, image patches can also be encoded as separate hypervectors [27]. As an example, imagine an image patch from the source domain that contains a car and pedestrian. We can assume without loss of generality that the image patch is represented as a hypervector:

$$v_{src} = c \otimes c_{src} + p \otimes p_{src}, \tag{1}$$

where $c$ and $p$ are hypervectors representing the car and pedestrian bound with source-specific hypervectors $c_{src}$ and $p_{src}$, respectively.

Since we do not have access to the underlying symbols associated with source and target domains, instead we must learn a mapping that unbinds source-specific information and binds target-specific information (and vice versa). Specifically, we would like to learn a hypervector mapping that can unbind these source-specific attributes and bind target-specific attributes in order to obtain a target representation $v_{tgt}$ as shown in Equation 2.

$$
\begin{aligned}
v_{src} \otimes u_{src \leftrightarrow tgt} &\approx v_{tgt}, \\
\text{where } u_{src \leftrightarrow tgt} &= (c_{src} \otimes c_{tgt} + p_{src} \otimes p_{tgt})
\end{aligned}
\tag{2}
$$

Since the binding/unbinding operation is invertible, $c_{src}$ and $p_{src}$ are unbound while $c_{tgt}$ and $p_{tgt}$ are correspondingly bound. Importantly, since these random hypervectors are very likely to be almost orthogonal, binding/unbinding the incorrect attributes (i.e., $c \otimes p_{src}$) simply adds noise to the resulting vector as demonstrated in Equation 3.

$$
\begin{aligned}
v_{src} \otimes u_{src \leftrightarrow tgt} &= c \otimes c_{src} \otimes (c_{src} \otimes c_{tgt} + p_{src} \otimes p_{tgt}) \\
&+ p \otimes p_{src} \otimes (c_{src} \otimes c_{tgt} + p_{src} \otimes p_{tgt}) \\
&= (c \otimes c_{tgt} + noise) + (noise + p \otimes p_{tgt}) \\
&= c \otimes c_{tgt} + p \otimes p_{tgt} + noise \approx v_{tgt}
\end{aligned}
\tag{3}
$$

It is also worth noting that $u_{src \leftrightarrow tgt}$ is equivalent to $u_{tgt \leftrightarrow src}$, which means a single vector can be used to map between the domains. Hence, $v_{tgt} \otimes u_{src \leftrightarrow tgt} \approx$
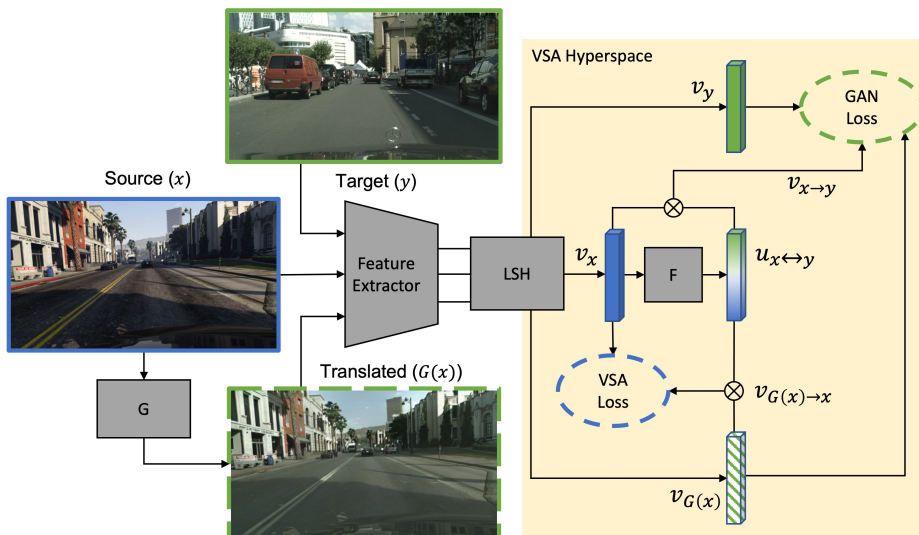
Fig. 3: Our method addresses semantic flipping by learning an invertible mapping in a high-dimensional vector space (VSA hyperspace) in order to ensure consistency between source and translated images. We extract features and use locality sensitive hashing (LSH) to encode source ($x$, blue), target ($y$, green), and translated ($G(x)$, striped green) images into this random VSA hyperspace. We use a GAN loss to train $G$ to generate images with hypervectors similar to those in the target domain. Finally, the hypervector mapping ($u_{x\leftrightarrow y}$, green-and-blue vertical bar) is used to invert our translation (Eqn 4) to recover the source hypervectors (VSA Loss). See Section 3.3 for more details.

$v_{src}$. Note that hypervectors in the VSA framework are distributed representations that are very robust to noise [1] and further that cosine similarity will still be high relative to similarity with random vectors. The above example demonstrates the unique properties of the VSA framework, which allow for algebraic formulations of symbolic representations.

## 3.2 VSA-Based Image Translation (VSAIT)

**Overview** Our goal is to translate images from a source to target domain, while minimizing semantic flipping caused by differences in the content distribution between the two domains. As shown in Figure 3, our method uses a hypervector adversarial loss (GAN Loss in Figure 3) that operates on VSA-based hypervector representations of image patches. In order to address semantic flipping we constrain the generator to preserve the content of the source domain via a VSA-based cyclic loss (VSA Loss in Figure 3). We do so by leveraging VSA to invert our translation in hypervector space to ensure consistency with the source content. We will cover these two loss functions later in this section. Our method has

three major network components– Generator $G$, Discriminator $D_Y$ and Mapper $F$ that we will discuss next. We present training details in Section 3.3.

**Generator** Our approach uses a generator $G : X \to Y$ (shown in Figure 3) and feature-level discriminator $D_Y$ in order to learn an unpaired image translation between $X$ (source) and $Y$ (target) domains, where $D_Y$ is trained to discriminate between hypervectors extracted from target and translated images.

As previously introduced in Section 3.1, to assign hypervectors to image patches, we follow Neubert et al. [28] and extract features using a pre-trained neural network (in our case, concatenated multiple layers of VGG19 [36]). These concatenated features are then randomly projected from the extracted feature vector $f_i \in \mathbb{R}^m$ to the random vector space $\mathbb{V} = [-1, 1]^n$ where $m >> n$. Using this approach, we denote hypervectors extracted from target images as $v_Y$, source images as $v_X$, and translated images as $v_{G(X)}$ (see Figure 3).

**Source $\leftrightarrow$ Target Mapper** Furthermore, we train a network $F$ to generate a hypervector mapping $u_{X \leftrightarrow Y}$ (e.g., Eqn 3) between source and target domains. We use this hypervector mapping to invert our translation, providing a VSA-based cyclic loss in hyperspace. The invertible mapping accomplishes two operations with a single step: unbinding of source (resp. target) representations and binding of target (resp. source) representations. This means that if we apply this mapping to our translated hypervectors $v_{G(X)}$, we should approximately obtain the original source hypervectors $v_X$. Similarly, if we apply this mapping to source hypervectors $v_X$, we approximately obtain translated hypervectors $v_{G(X)}$. We denote translated hypervectors that have been mapped to the source domain as $v_{G(X) \to X}$ and those mapped from source to target domain as $v_{X \to Y}$ as shown in Equation 4 as well as Figure 3.

$$
\begin{aligned}
v_{G(X)} \otimes u_{X \leftrightarrow Y} = v_{G(X) \to X} \approx v_X \\
v_X \otimes u_{X \leftrightarrow Y} = v_{X \to Y} \approx v_{G(X)}
\end{aligned}
\tag{4}
$$

**Hypervector Adversarial Loss** In order to train $G$, $D_Y$ and $F$ to translate images that match the target distribution, we use an adversarial loss [10]. Specifically, we want the hypervectors of our translated image $v_{G(X)}$ to be similar to those from the target domain $v_Y$. Furthermore, by binding the hypervector mapping $u_{X \leftrightarrow Y}$ to source vectors $v_X$, we should obtain a hypervector of the source features mapped to the target domain (Eqn 4) . Therefore, we train $F$ along with $G$ and $D_Y$ using the following adversarial loss:

$$
\begin{aligned}
\mathcal{L}_{GAN}(G, D_Y, F, X, Y) = &\mathbb{E}_{y \sim p_Y(y)}[\log D_Y(v_y)] \\
&+ \mathbb{E}_{x \sim p_X(x)}[\log(1 - D_Y(v_{G(x)}))] \\
&+ \mathbb{E}_{x \sim p_X(x)}[\log(1 - D_Y(v_x \otimes F(v_x)))],
\end{aligned}
\tag{5}
$$

where $G$ and $F$ networks are trained to fool the discriminator $D_Y$, thereby minimizing the objective while $D_Y$ attempts to maximize it.

**VSA-Based Cyclic Loss** Although adversarial training provides a method to generate images that match the target distribution, the differences in content between domains will result in semantic flipping. We therefore need a loss that constrains the generator to preserve source content and reduce semantic flipping. To do so, we incorporate our VSA-based cyclic loss, which ensures that the same hypervectors are obtained when mapping translated vectors back from target to source domain ($v_{G(X) \to X}$):

$$\mathcal{L}_{VSA}(G, X) = \mathbb{E}_{x \sim p_X(x)} \left[ \frac{1}{n} \sum_{i=1}^{n} dist \left( v_x^i, v_{G(x) \to x}^i \right) \right],      \tag{6}$$

where $dist(\cdot, \cdot)$ is the cosine distance averaged across the $n$ hypervectors (representing image patches indexed by $i$) for the image $x$. By minimizing the cosine distance ($1 -$ cosine similarity) between $v_{G(X) \to X}$ and $v_X$, we ensure that the same features are recovered after inverting our translation (i.e., representations of "car" and "pedestrian" are maintained in the example from Equation 1).

Our overall objective combines our adversarial and VSA-based cyclic losses, training $G$ to generate images matching the target domain and $F$ to invert the translation to ensure consistency with the source domain:

$$\mathcal{L}(G, D_Y, F, X, Y) = \mathcal{L}_{GAN}(G, D_Y, F, X, Y) + \lambda \mathcal{L}_{VSA}(G, X),      \tag{7}$$

where $\lambda$ controls the relative importance of our VSA-based cyclic loss.

### 3.3   VSAIT Training

We train our method with a GAN-based training framework using the objective described in Equation 7. We begin by randomly sampling an image from each domain and translating the source image into the target domain via the generator $G$. In order to use the VSA framework, we must first encode data into the hypervector space. Therefore, we extract features from each image (source, target, and translated) at multiple layers of a pre-trained neural network (Feature Extractor in Figure 3) as the first step in encoding image patches into the VSA hyperspace. We consider each image patch to be represented by the extracted features sharing its receptive field. We flatten and concatenate these features across layers, resulting in a set of feature vectors (one per image patch) with a dimensionality of $m$.

We then use LSH (Figure 3) to reduce the dimensionality of the extracted feature vectors into our random hyperspace $\mathbb{V}$, which reduces the dimensionality to $n << m$. Specifically, we normalize and project the extracted feature vectors using a random matrix where each row is sampled from a standard normal distribution and normalized to unit length. The resulting vectors are therefore in the range $[-1, 1]$, which is necessary to implement the VSA binding operation as described in Section 3.1. These hypervectors represent the image patches for the source, target, and translated images as shown in Figure 3. Note that this process does not require training nor is it necessarily dependent on a specific feature extractor, although different feature extractors may encode different information.

Finally, we use $F$ to generate a hypervector mapping for each source hypervector. We use these hypervector mappings ($u_{x\leftrightarrow y}$ in Figure 3) to unbind source (resp. target) information and correspondingly bind target (resp. source) information (Eqn 4). By applying this mapping to the source hypervectors, we should obtain a hypervector representation of the source image translated into the target domain ($v_{x\rightarrow y}$ in Figure 3), which is used in our adversarial loss (Eqn 5) to train $F$. The ultimate goal of this step, however, is to use the mapping to invert our translation $v_{G(x)}$ and ensure that we recover the same source information $v_x$ (Eqn 6). Doing so will reduce semantic flipping by constraining $G$ to generate images that have hypervectors that are consistent with the source domain.

## 4   Experiments

We evaluate our method using experiments across multiple datasets. First, we compare against baseline techniques for GTA [34] to Cityscapes [6], where the two domains have inherent semantic differences. We then perform experiments using datasets sub-sampled to create differences in semantic statistics. Specifically, we follow the method in [14] to sub-sample the Google Maps [13] dataset, which is typically designed for paired image translation tasks. We show that VSAIT works as intended to reduce semantic flipping while still generating diverse image translations. Our qualitative results in Figure 2 demonstrate that the hyperspace mapping constrains the generator to translate features that are consistent between source and target domains. Furthermore, we show that other methods still have significant artifacts and hallucinations related to semantic flipping (Figure 4). Our quantitative results show improvements against the other baselines, particularly for the GTA to Cityscapes task. Finally, we perform an ablation study to evaluate the contributions of different components in VSAIT.

### 4.1   Implementation Details

We follow CUT [30] in our choice of the generator network architecture. For the discriminator network, we use a three-layer fully-convolutional network with $1 \times 1$ convolutional filters. For the mapping network $F$, we use a two-layer fully-convolutional network. We train VSAIT using the Adam optimizer [19] ($\beta_1 = 0$, $\beta_2 = 0.9$) with a batch size of 1 and learning rate of 0.0001 for the generator (and mapping network $F$) and 0.0004 for the discriminator. To improve adversarial training, we use the "hinge loss" [24] rather than the negative log-likelihood objective in $\mathcal{L}_{GAN}$ (Eqn 5). For GTA to Cityscapes, we use $\lambda = 10$ in Equation 7 and $\lambda = 5$ for other experiments. See Supplemental for more details.

### 4.2   Datasets

We perform our experiments using three datasets: GTA [34], Cityscapes [6], and Google Maps [13]. GTA [34] is a synthetic dataset of 24966 images generated from the video game Grand Theft Auto V. Cityscapes [6] is a real dataset of

driving scenarios accompanied with finely-annotated semantic segmentation labels, which includes 2975 training and 500 validation images. The Google Maps dataset [13] is a collection of 2194 paired maps and aerial photos (1096 training and 1098 validation images).

### 4.3    Baselines

We compare our method against several baselines that are relevant to the semantic flipping problem: GcGAN, DRIT, CUT, SRUNIT, and EPE. We choose these methods for comparison as they reflect the recent approaches focused on reducing semantic flipping. Specifically, they represent methods that use image-level consistency (GcGAN [8]), domain-invariant encoding (DRIT [23], CUT [30], SRUNIT [14]), and task-level consistency (EPE [33]). It is worth reiterating that EPE only works for synthetic datasets with access to G-buffers (e.g., albedo, glossiness, depth, normals), which are not publicly available for GTA and do not exist for real datasets (e.g., Google Maps). Therefore we compare qualitatively to EPE and rely on quantitative values reported in their work.

### 4.4    GTA $\rightarrow$ Cityscapes

We first demonstrate our method using GTA to Cityscapes, which are unpaired datasets that naturally differ in their semantic statistics. For training, we use 20000 of GTA images for our source dataset and all 2975 training images for our target dataset. Following [33], we evaluate our image translation performance using the Kernel Inception Distance (KID) [4] which measures the distance of extracted features from translated and target images using the pre-trained InceptionV3 network [37]. As seen in the Table 2, our method outperforms the baseline methods in the KID metric. As shown in Figure 4, only EPE has similar quality of image translations for GTA to Cityscapes. However, whereas other methods hallucinate trees and Mercedes hood ornaments, EPE has a tendency to remove palm trees and add unnatural lighting to cars even in dark scenes.

We additionally evaluate translation quality via semantic segmentation metrics computed using DeepLab V3 [5] pretrained on Cityscapes. In line with [14], we report three metrics related to semantic segmentation performance (Table 1): pixel accuracy, class accuracy, and mean intersection over union (mIoU). As shown in Table 1, we outperform the other baselines by a wide margin. Whereas KID computes the distance between translated and target images in feature space, semantic segmentation metrics better reflect semantic flipping directly. As seen in Figure 4, other methods often generate trees and other features in the sky, which will lower semantic segmentation performance in those regions.

### 4.5    Google Map $\rightarrow$ Aerial Photo

We then demonstrate performance on the Google Maps [13] dataset. Here we sub-sample the 1096 training images using K-means clustering of the histograms

Fig. 4: We compare our VSAIT translation to those of baseline methods on an example GTA image. Typical methods often suffer from semantic flipping in open sky regions of GTA, as such regions are observed less often in Cityscapes (i.e., content gap). Here we see that SRUNIT [14], CUT [30], GcGAN [8] and DRIT [23] each suffer from sky hallucinations, while EPE [33] removes palm trees (features that are absent from Cityscapes) in the distance. Meanwhile, our method is able to translate to Cityscapes style while preserving semantics from the original GTA image.

from grayscale map images to obtain two datasets with different semantic statistics (as in [14]). We evaluate translation performance on all 1098 validation images and report three pixel-level metrics (Table 1). The first metric is the L2 distance between translated and target images. We also report pixel accuracy defined as the percentage of pixels with a maximum absolute difference less than a given threshold (i.e., $max(|r_i - r_i'|, |g_i - g_i'|, |b_i - b_i'|) < \delta$). For the map to aerial

Table 1: Quantitative evaluation across datasets with unmatched semantics. The metrics included are average pixel prediction accuracy (pxAcc), average class prediction accuracy (clsAcc), mean IoU (mIoU), average L2 distance (Dist), and pixel accuracy with task-specific thresholds (Acc).

| Method | GTA → Cityscapes | | | Map → Photo | | | Photo → Map | | |
|---|---|---|---|---|---|---|---|---|---|
| | pxAcc | clsAcc | mIoU | Dist | $Acc(\delta_1)$ | $Acc(\delta_2)$ | Dist | $Acc(\delta_1)$ | $Acc(\delta_2)$ |
| GcGAN [8] | 65.62 | 32.38 | 22.64 | 71.47 | 28.87 | 43.48 | 23.62 | 15.00 | 30.65 |
| DRIT [23] | 64.28 | 32.17 | 20.99 | 70.87 | 28.97 | 43.56 | 24.19 | 13.94 | 29.01 |
| CUT [30] | 64.59 | 32.19 | 20.35 | 70.28 | 28.86 | 44.07 | 23.44 | 16.25 | 31.34 |
| SRUNIT [14] | 67.21 | 32.97 | 22.69 | 68.55 | 30.41 | 45.91 | 23.00 | **17.67** | **32.78** |
| VSAIT (ours) | **76.48** | **45.33** | **30.89** | **64.98** | **45.3** | **69.28** | **22.86** | 15.2 | 32.13 |

photo task, we use thresholds of $\delta = 30, 50$ [14]. Similar to GTA to Cityscapes, our approach substantially outperforms the baseline methods. Particularly interesting is the improvement we see in the pixel accuracy for both thresholds ($\delta = 30, 50$), whereas the other methods perform very similarly.

## 4.6   Aerial Photo → Google Map

We additionally demonstrate performance for the task of photo to map using the same sub-sampled datasets obtained from the Google Maps training dataset as described above. For this experiment, we evaluate the same metrics as with map to photo but with pixel accuracy threshold of $\delta = 3, 5$ [14] (Table 1). Although we outperform the baseline methods on the L2 distance metric (Dist in Table 1), we did not observe improvements in the pixel accuracy metrics. We do believe that our approach using VSA is capable of improving these metrics with a more suitable encoding method (as opposed to VGG) for images that have regions without contours or complex features (as is the case for some regions in Google Map images representing landscape or water). However, we leave study of encoding methods for VSA in computer vision applications to future research.

## 4.7   Ablation Study

We demonstrate the effect of VSAIT by evaluating the contribution of the VSA-based cyclic loss (Eqn 6), the learned hypervector mapping $u_{X \leftrightarrow Y}$, and the hypervector dimensionality. We use the GTA to Cityscapes task to evaluate these ablations on semantic segmentation performance metrics (Table 2). We first demonstrate that by removing the VSA-based cyclic loss, the generator learns to translate images without preserving content, demonstrating the serious problem of semantic flipping (Figure 5B). Next, we show the importance of the invertible mapping $u_{X \leftrightarrow Y}$ by generating a random hypervector mapping instead of learning the mapping adversarially (Eqn 5). When using a random hypervector instead of the learned mapping generated by $F$, the generator maintains global structure but since the invertible mapping does not recover the source hypervector (Eqn

Table 2: Evaluations on GTA to Cityscapes. Left: quantitative comparison for KID metric. *Mean values reported in [33]. Right: ablation study task demonstrating the contributions for each component of VSAIT method.

| Method | KID |
|---|---|
| CUT | 21.18* |
| SRUNIT | 16.27 |
| DRIT | 14.69 |
| GcGAN | 12.32 |
| EPE | 10.95* |
| VSAIT (ours) | **8.74** |

| Ablation | pxAcc | clsAcc | mIoU |
|---|---|---|---|
| w/o VSA Loss | 52.17 | 16.14 | 10.21 |
| Random Hypervector Mapping | 65.75 | 29.32 | 17.53 |
| Reduced Hypervector Dim | 66.42 | 40.34 | 25.04 |
| VSAIT (ours) | **76.48** | **45.33** | **30.89** |

4) the local content is often changed (semantic flipping shown in Figure 5C). Finally, we demonstrate the importance of using the VSA hypervector space by reducing the dimensionality of $\mathbb{V}$ from 4096 to 128. Reducing the dimensionality of the hyperspace from 4096 to 128 results in noisy image translations that seem to only reflect global changes in style. However, even using a low-dimensional hypervector VSAIT still outperforms most methods compared in Table 1.

## 5   Conclusion

In this paper, we address the semantic flipping problem in unpaired image translation using a novel approach based on VSA framework [9, 15]. We show impor-
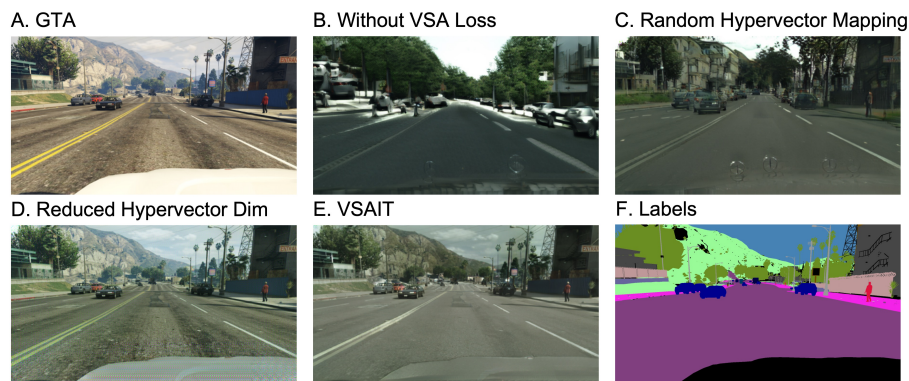


Fig. 5: Ablation study for GTA to Cityscapes. (A) and (F) show the ground truth image and semantic labels for the GTA source image, respectively. (B) demonstrates the challenge of semantic flipping when using GAN-based methods without constraining the generator. (C) shows the importance of learning the hypervector mapping. (D) demonstrates the importance of high-dimensional space for this method. (E) shows the effect on image translation when incorporating all components of our approach.

tant qualitative and quantitative improvements over previous methods that have attempted to address this problem, demonstrating that VSA can be used to invert image translations and ensure consistency with the source domain. Given its inherent versatility, we hope this work inspires future research using VSA in more computer vision applications.

# References

[1]   Subutai Ahmad and Jeff Hawkins. "Properties of sparse distributed representations and their application to hierarchical temporal memory". In: *arXiv preprint arXiv:1503.07469* (2015).

[2]   Amjad Almahairi et al. "Augmented cyclegan: Learning many-to-many mappings from unpaired data". In: *International Conference on Machine Learning.* PMLR. 2018, pp. 195–204.

[3]   Sagie Benaim and Lior Wolf. "One-sided unsupervised domain mapping". In: *Advances in neural information processing systems* 30 (2017).

[4]   Mikołaj Bińkowski et al. "Demystifying mmd gans". In: *arXiv preprint arXiv:1801.01401* (2018).

[5]   Liang-Chieh Chen et al. "Rethinking atrous convolution for semantic image segmentation". In: *arXiv preprint arXiv:1706.05587* (2017).

[6]   Marius Cordts et al. "The Cityscapes Dataset for Semantic Urban Scene Understanding". In: *CVPR.* June 2016.

[7]   Jeevan Devaranjan, Amlan Kar, and Sanja Fidler. "Meta-sim2: Unsupervised learning of scene structure for synthetic data generation". In: *European Conference on Computer Vision.* Springer. 2020, pp. 715–733.

[8]   Huan Fu et al. "Geometry-consistent generative adversarial networks for one-sided unsupervised domain mapping". In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 2019, pp. 2427–2436.

[9]   Ross W Gayler. "Vector symbolic architectures answer Jackendoff's challenges for cognitive neuroscience". In: *arXiv preprint cs/0412059* (2004).

[10]  Ian Goodfellow et al. "Generative adversarial nets". In: *Advances in neural information processing systems* 27 (2014).

[11]  Judy Hoffman et al. "Cycada: Cycle-consistent adversarial domain adaptation". In: *International conference on machine learning.* PMLR. 2018, pp. 1989–1998.

[12]  Xun Huang et al. "Multimodal unsupervised image-to-image translation". In: *Proceedings of the European conference on computer vision (ECCV).* 2018, pp. 172–189.

[13]  Phillip Isola et al. "Image-to-image translation with conditional adversarial networks". In: *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2017, pp. 1125–1134.

[14]    Zhiwei Jia et al. "Semantically Robust Unpaired Image Translation for Data with Unmatched Semantics Statistics". In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 14273–14283.

[15]    Pentti Kanerva. "Hyperdimensional computing: An introduction to computing in distributed representation with high-dimensional random vectors". In: *Cognitive computation* 1.2 (2009), pp. 139–159.

[16]    Amlan Kar et al. "Meta-sim: Learning to generate synthetic datasets". In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2019, pp. 4551–4560.

[17]    SJ Kent and BA Olshausen. "A vector symbolic approach to scene transformation". In: *Cognitive computational neuroscience (ccn'17)(extended abstract)[link]* (2017).

[18]    Taeksoo Kim et al. "Learning to discover cross-domain relations with generative adversarial networks". In: *International conference on machine learning*. PMLR. 2017, pp. 1857–1865.

[19]    Diederik P Kingma and Jimmy Ba. "Adam: A Method for Stochastic Optimization". In: *ICLR*. 2015.

[20]    Denis Kleyko et al. "A Survey on Hyperdimensional Computing aka Vector Symbolic Architectures, Part I: Models and Data Transformations". In: *arXiv preprint arXiv:2111.06077* (2021).

[21]    Denis Kleyko et al. "A Survey on Hyperdimensional Computing aka Vector Symbolic Architectures, Part II: Applications, Cognitive Models, and Challenges". In: *arXiv preprint arXiv:2112.15424* (2021).

[22]    Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. "Imagenet classification with deep convolutional neural networks". In: *Advances in neural information processing systems* 25 (2012).

[23]    Hsin-Ying Lee et al. "Diverse image-to-image translation via disentangled representations". In: *Proceedings of the European conference on computer vision (ECCV)*. 2018, pp. 35–51.

[24]    Jae Hyun Lim and Jong Chul Ye. "Geometric gan". In: *arXiv preprint arXiv:1705.02894* (2017).

[25]    Ming-Yu Liu, Thomas Breuel, and Jan Kautz. "Unsupervised image-to-image translation networks". In: *Advances in neural information processing systems* 30 (2017).

[26]    Guglielmo Montone, J Kevin O'Regan, and Alexander V Terekhov. "Hyperdimensional computing for a visual question-answering system that is trainable end-to-end". In: *arXiv preprint arXiv:1711.10185* (2017).

[27]    Peer Neubert and Stefan Schubert. "Hyperdimensional computing as a framework for systematic aggregation of image descriptors". In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2021, pp. 16938–16947.

[28]    Peer Neubert, Stefan Schubert, and Peter Protzel. "An introduction to hyperdimensional computing for robotics". In: *KI-Künstliche Intelligenz* 33.4 (2019), pp. 319–330.

[29]    Evgeny Osipov et al. "HyperSeed: Unsupervised Learning with Vector Symbolic Architectures". In: *arXiv preprint arXiv:2110.08343* (2021).

[30]    Taesung Park et al. "Contrastive learning for unpaired image-to-image translation". In: *European Conference on Computer Vision*. Springer. 2020, pp. 319–345.

[31]    Aayush Prakash et al. "Self-Supervised Real-to-Sim Scene Generation". In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 16044–16054.

[32]    Scott Purdy. "Encoding data for HTM systems". In: *arXiv preprint arXiv:1602.05925* (2016).

[33]    Stephan R Richter, Hassan Abu AlHaija, and Vladlen Koltun. "Enhancing Photorealism Enhancement". In: *arXiv preprint arXiv:2105.04619* (2021).

[34]    Stephan R Richter et al. "Playing for data: Ground truth from computer games". In: *European conference on computer vision*. Springer. 2016, pp. 102–118.

[35]    Kenny Schlegel, Peer Neubert, and Peter Protzel. "A comparison of vector symbolic architectures". In: *Artificial Intelligence Review* (2021), pp. 1–33.

[36]    Karen Simonyan and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition". In: *arXiv preprint arXiv:1409.1556* (2014).

[37]    Christian Szegedy et al. "Rethinking the inception architecture for computer vision". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 2818–2826.

[38]    Zili Yi et al. "Dualgan: Unsupervised dual learning for image-to-image translation". In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 2849–2857.

[39]    Jun-Yan Zhu et al. "Unpaired image-to-image translation using cycle-consistent adversarial networks". In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 2223–2232.

## A    Additional Implementation Details

In order to encode image patches into the VSA hyperspace, we extract features using VGG19 [5] pre-trained on ImageNet [1]. We select multiple layers from VGG19 and concatenate features that share the same receptive field. For all experiments, we extract features from multiple convolutional layers in VGG19. For GTA to Cityscapes, we concatenate features for patch sizes $16 \times 16$, $8 \times 8$, $4 \times 4$, $2 \times 2$, and, $1 \times 1$. This means that extracted features within each patch are concatenated into a single feature vector, which is reduced in dimensionality using locality sensitive hashing as described in Section 3.3 of the main paper. For all experiments, we used a hypervector dimensionality of 4096, which is the dimensionality of the input to the discriminator and mapping networks. We then used 1024 channels in all layers except the last (1 channel) for the discriminator, and 4096 channels for both layers in the mapping network.

For Google Map to Aerial Photo, we use the same patch sizes as for GTA to Cityscapes. However, for Aerial Photo to Google Map, we empirically found that using "dilated" patches gave better performance. For this experiment, we use features extracted from four convolutional layers with patch sizes $32 \times 32$, $16 \times 16$, $8 \times 8$, and $4 \times 4$, where each patch has spaces between locations in the extracted feature maps (i.e., equivalent to a convolutional dilation factor of 3).

For all experiments, we trained using $256 \times 256$ images (as was done in the baseline methods). In the case of GTA to Cityscapes, we first resize the images to a height and width of $256 \times 512$ followed by random cropping to $256 \times 256$. For the Google Map and Aerial Photo datasets, images are resized from $600 \times 600$ to $256 \times 256$. Note that we do not report results for EPE in Table 1 of the main paper as we are not able to replicate their method using $256 \times 256$ image resolution for GTA to Cityscapes and it is not feasible to use their method for datasets without G-buffers. As was done in [3], we train for 20 epochs for GTA to Cityscapes and 400 epochs for Google Map to Aerial Photo and Aerial Photo to Google Map experiments. We multiply the learning rates (defined in Section 4.1 of the main paper) by a factor of 0.5 every 100K iterations.

## B    Difference between VSA-Based Cyclic Loss and Perceptual Loss

Many image-to-image translation methods have included a so-called "perceptual loss" to minimize the distance between translated and source images in VGG feature space. Although our VSA-based cyclic loss (Eqn 6) may seem similar to the perceptual loss, our method is distinct: we minimize the distance between source hypervectors and translated hypervectors that are mapped back from the target to source domain, whereas methods using a perceptual loss directly minimize the distance between source and translated features. Furthermore, as noted previously, encoding these source and translated features into the VSA hyperspace gives greater assurance that different semantic content will be almost orthogonal to each other.

## C  Additional Ablation Studies

As indicated above, we changed the patch sizes used during feature extraction across different experiments (i.e., "dilated" patches in Aerial Photo to Google Map). We therefore, also looked at the effect of patch size within a specific experiment. In addition to the $16 \times 16$ patch size used in the GTA to Cityscapes experiment, we also tested $32 \times 32$ and $8 \times 8$ patch sizes. Compared to the results in Table 1 (mIoU 30.89), we observed mIoU of 32.29 and 31.47 for 32 and 8 patch sizes, respectively.

Furthermore, we conducted an ablation of the weight ($\lambda$) for our VSA-based cyclic loss, showing that $\lambda < 5$ results in semantic flipping, which is reflected in lower mIoU for the GTA to Cityscapes experiment. Specifically, we observed lower mIoU for smaller (21.55 and 26.49 for $\lambda = 1, 2$) compared to larger $\lambda$ (29.13, 30.89, 29.78 for $\lambda = 5, 10, 20$).

Finally, we also tested the effect of the hypervector dimensionality. As demonstrated in [4], the hypervector dimensionality should be at least greater than 700 (the probability of two randomly sampled vectors being orthogonal approaches 1). We tested 512 and 2048 dimensional hypervectors, demonstrating worse performance for 512 compared to 2048 (29.35 vs. 30.87 mIoU for GTA to Cityscapes, respectively).

## D  Hypervector Adversarial Loss

Since our hypervector adversarial loss (Eqn 5) has two negative targets and one positive target, it's possible that this could lead to an imbalance in training. For the current experiments, we have tested using a balanced weighting of the loss terms and did not empirically observe differences in stability or performance. However, this may be an important consideration for future research with other datasets.

## E  Quantitative Comparison to Additional Methods for GTA to Cityscapes

In order to compare our method against other relevant methods that have also been applied to GTA to Cityscapes, we trained semantic segmentation networks using our translated GTA images as done in other methods. In these comparisons we used ResNet-50 as the feature extractor rather than VGG19. Specifically, we compared our method against CyCADA [2] and Fourier Domain Adaptation (FDA) [7]. CyCADA [2] is a method that uses the cycle-consistency approach from CycleGAN [8] in order to ensure semantic consistency between segmentation predictions before and and after image translation. Unlike GAN-based methods, FDA [7] computes the Fourier Transform of source and target images, and replaces the amplitudes of low-frequency features in the source images with those from the target images. We observed comparable performance to these approaches (45.7 mIoU for VSAIT vs. 42.7 in CyCADA [2] and 44.6 in FDA [7]).

We further compared our results to those from [6], which is a recent video-to-video translation technique that aims to prevent semantic flipping via pseudo-supervision with optical flow from sequential video frames. By using a checkpoint from the authors to evaluate semantic segmentation performance (GTA to Cityscapes, Table 1), we observed lower performance compared to our method (pixel accuracy: 66.04, class accuracy: 38.42, and mIoU: 24.53).

## F   VSA Hypervector Encoding

As noted in the main paper, VGG19 works well for most natural image domains; however, it may not be optimal for encoding semantic segmentation masks or similar types of images (e.g., Google Map images). For example, using VSAIT to translate semantic masks to Cityscapes images using the same sub-sampling method used in [3], we observed worse performance compared to that reported in [3] (pixel accuracy: 69.28, class accuracy: 31.37, mIoU: 21.59). Therefore, an important future direction for this research is to understand and improve encoding methods when using VSA for computer vision applications.

## G   Additional Visual Examples

### G.1   GTA → Cityscapes

In addition to the comparisons shown in the main paper, we provide more qualitative examples showing our method relative to the baseline methods for GTA to Cityscapes. As described in the main paper, GTA and Cityscapes have semantics that are naturally different. Most notably, Cityscapes has more vegetation in the upper half of the image whereas GTA has more sky (see also Figure 1 from the main paper). Therefore, many GAN-based image translation methods hallucinate trees in the sky as shown in Figures 6 and 7. More interestingly, EPE does not hallucinate trees but instead removes palm trees from images. This is likely due to the fact that palm trees do not appear in Cityscapes. On the other hand, our method correctly retains the palm trees without hallucinating more trees in the sky.

Figure 8 provides a similar example of the typical hallucinations that occur when translating GTA images with broad sky regions. However, notice the sandy hill on the left side of the image, which we highlight in the figure. As shown in the top-right panel, the semantic segmentation label for this region is the same for both the grassy and sandy portions of the hill. Since EPE uses these labels as part of its task-level consistency loss, it generates the same feature (i.e., grass) across the region, irrespective of the underlying specific content differences.

Fig. 6: Examples of semantic flipping for baseline methods in GTA to Cityscapes experiment.
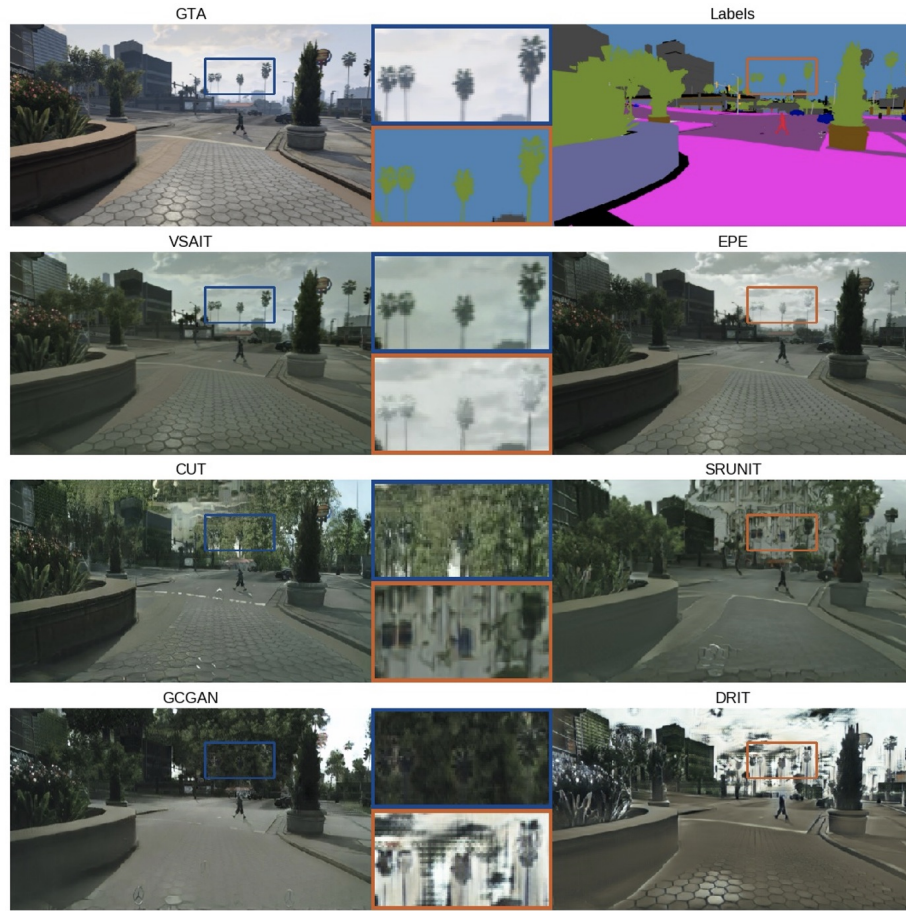
Fig. 7: Examples of semantic flipping for baseline methods in GTA to Cityscapes experiment.
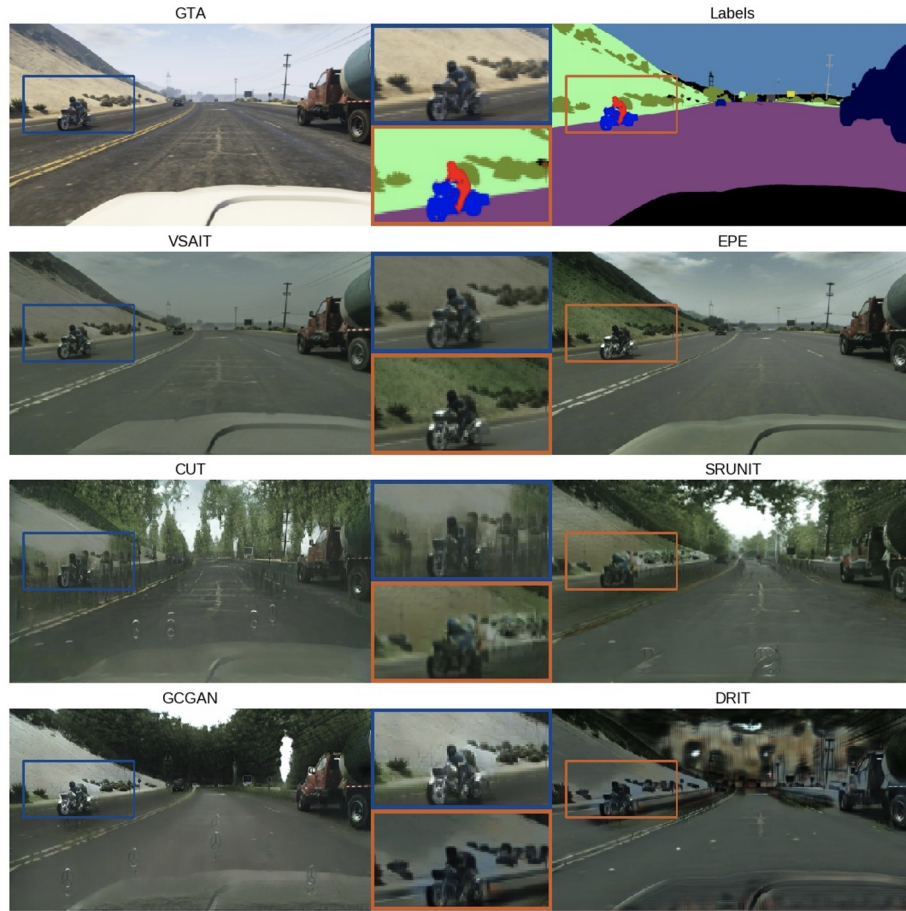
Fig. 8: Examples of semantic flipping for baseline methods in GTA to Cityscapes experiment.
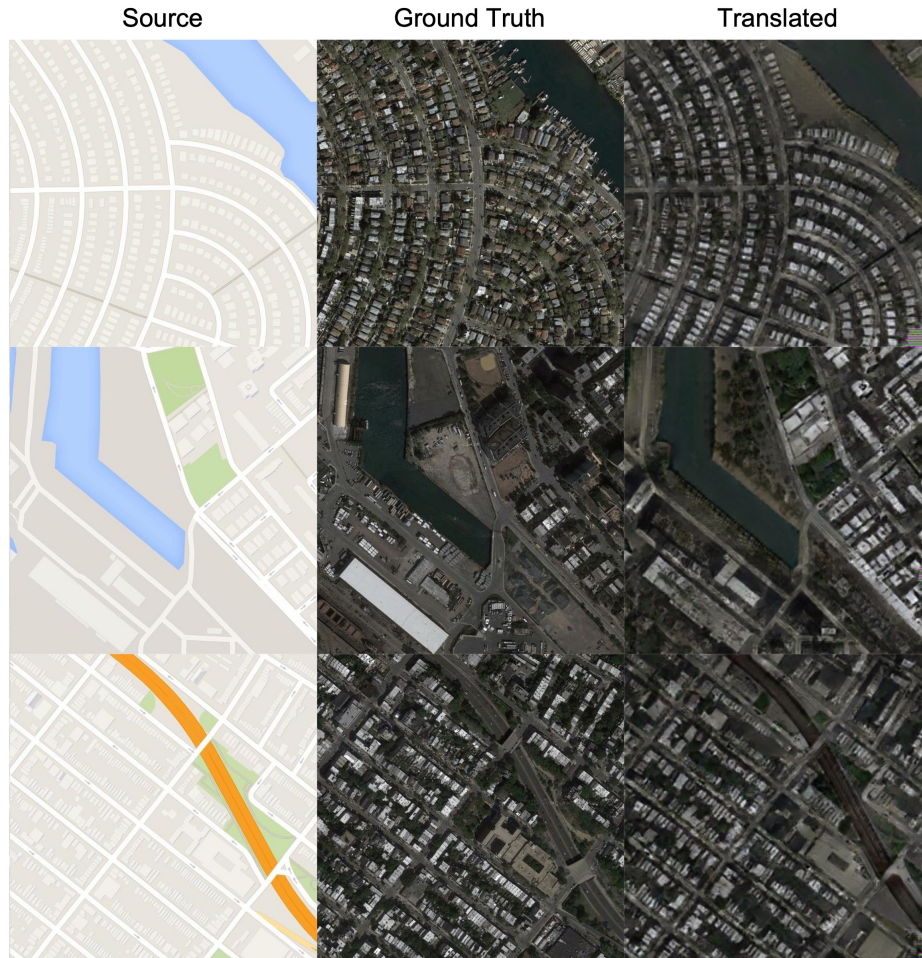
## G.2    Google Map → Aerial Photo



Fig. 9: Example image translations using VSAIT for the Google Map to Aerial Photo experiment.

As described in Section 4.5 of the main paper, we sub-sample the Google Map and Aerial Photo datasets (following [3]) to obtain unpaired source and target datasets that differ in their semantic statistics. As a result, the sub-sampled datasets differ in the proportion of image area containing land versus water. Therefore, similar to the GTA to Cityscapes experiment, we would expect hallucinations of houses or trees in the water. We demonstrate across multiple examples in Figures 9 and 10 that our method does not exhibit semantic flipping. It's

possible that some methods may rely on a one-to-one mapping using the RGB values in the map images (e.g., blue for water) that could reduce the diversity of translations overall. However, as shown in Figure 11, we observe diverse translations for the same RGB values across different map images. Specifically, notice the translations associated with gray pixels in each source image as highlighted in the figure.
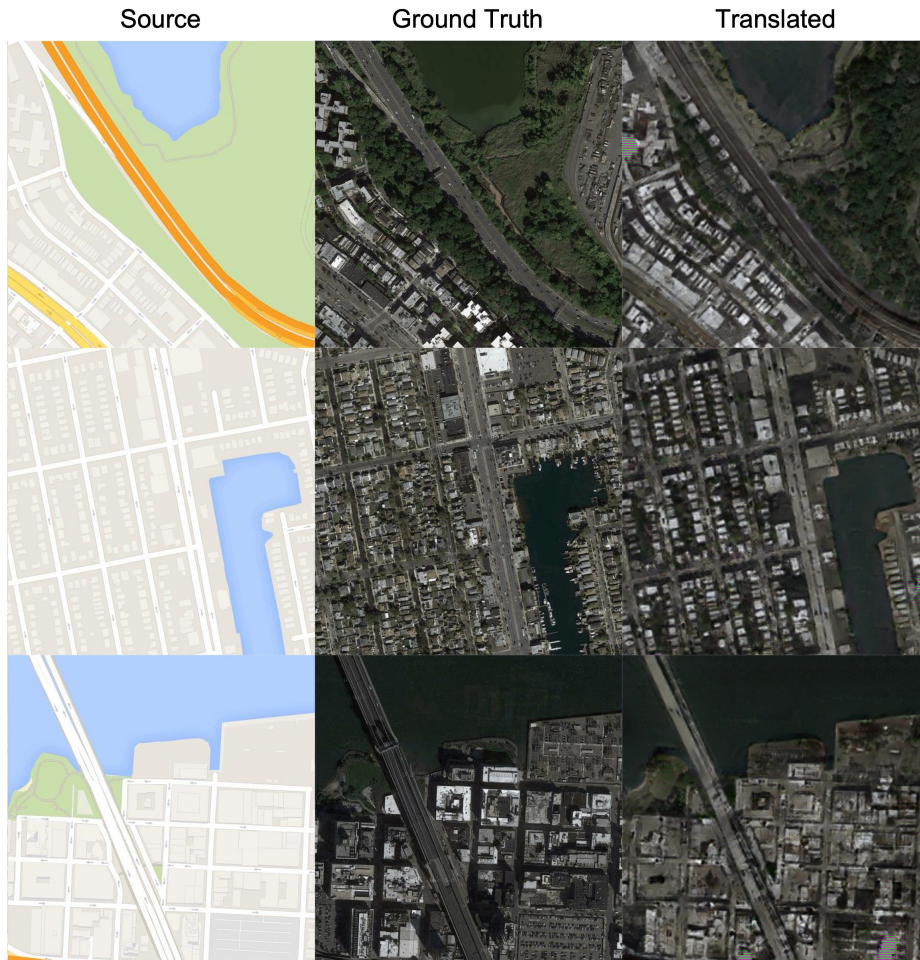


Fig. 10: Example image translations using VSAIT for the Google Map to Aerial Photo experiment.

Source              Ground Truth              Translated



Fig. 11: Example image translations using VSAIT for the Google Map to Aerial Photo experiment.

### G.3    Aerial Photo → Google Map

Finally, we demonstrate examples of our image translations from the Aerial Photo to Google Map experiment. As mentioned previously, we use the same sub-sampling procedure to obtain datasets with unmatched semantic statistics. Since the sub-sampled Google Map dataset has more regions with water relative to the sub-sampled Aerial Photo dataset, we would expect that GAN-based methods might suffer from semantic flipping by hallucinating water. However, as shown in Figures 12 and 13, our method learns the correct mapping between water and does not exhibit semantic flipping in scenes without water (Figure 14).
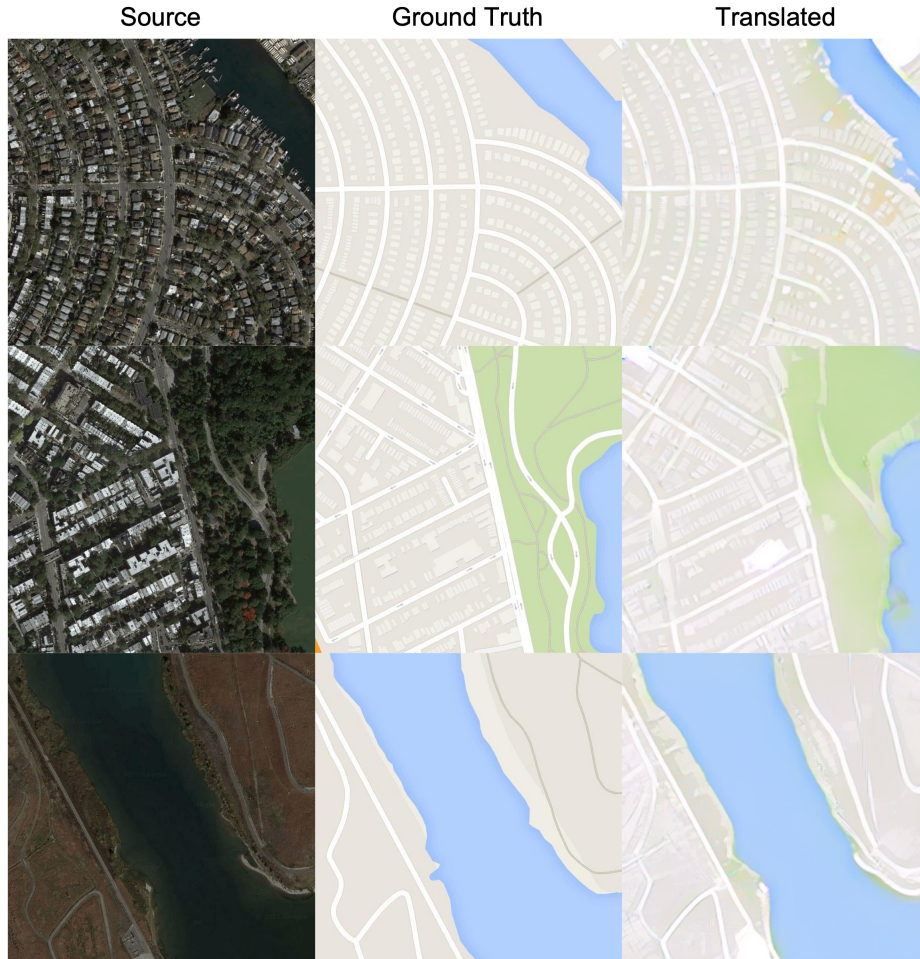


Fig. 12: Example image translations using VSAIT for the Aerial Photo to Google Map experiment.

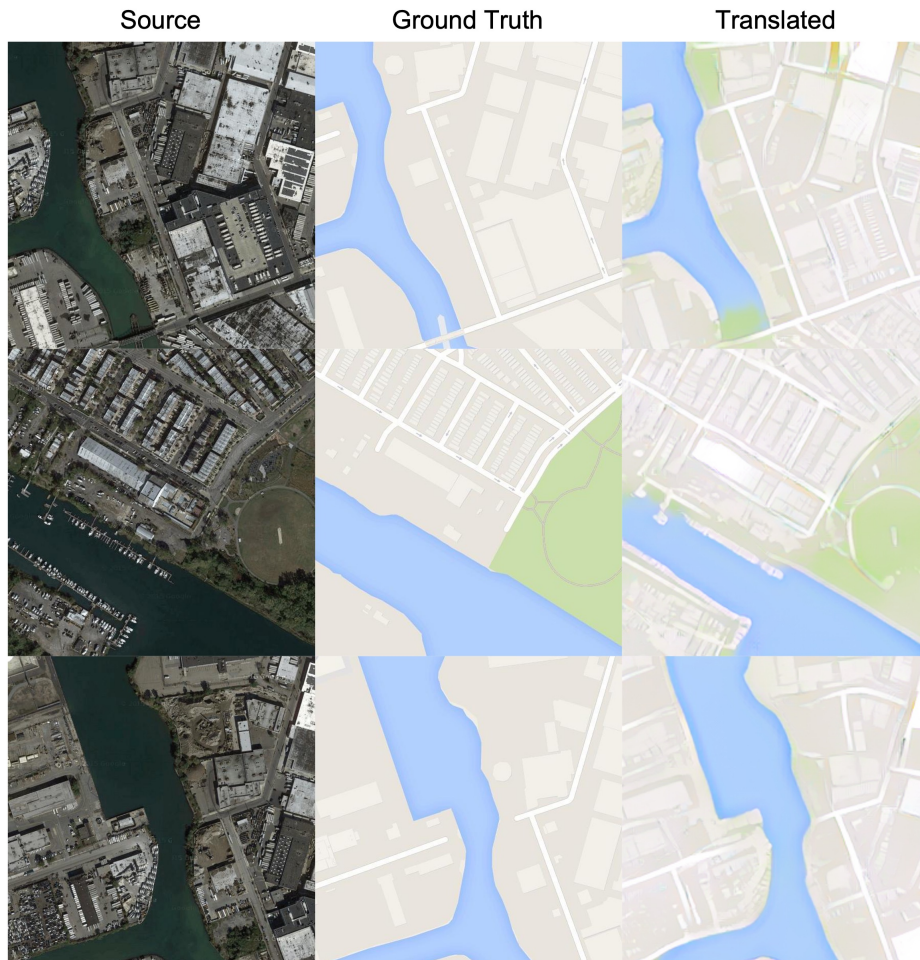| Source | Ground Truth | Translated |
|--------|--------------|------------|

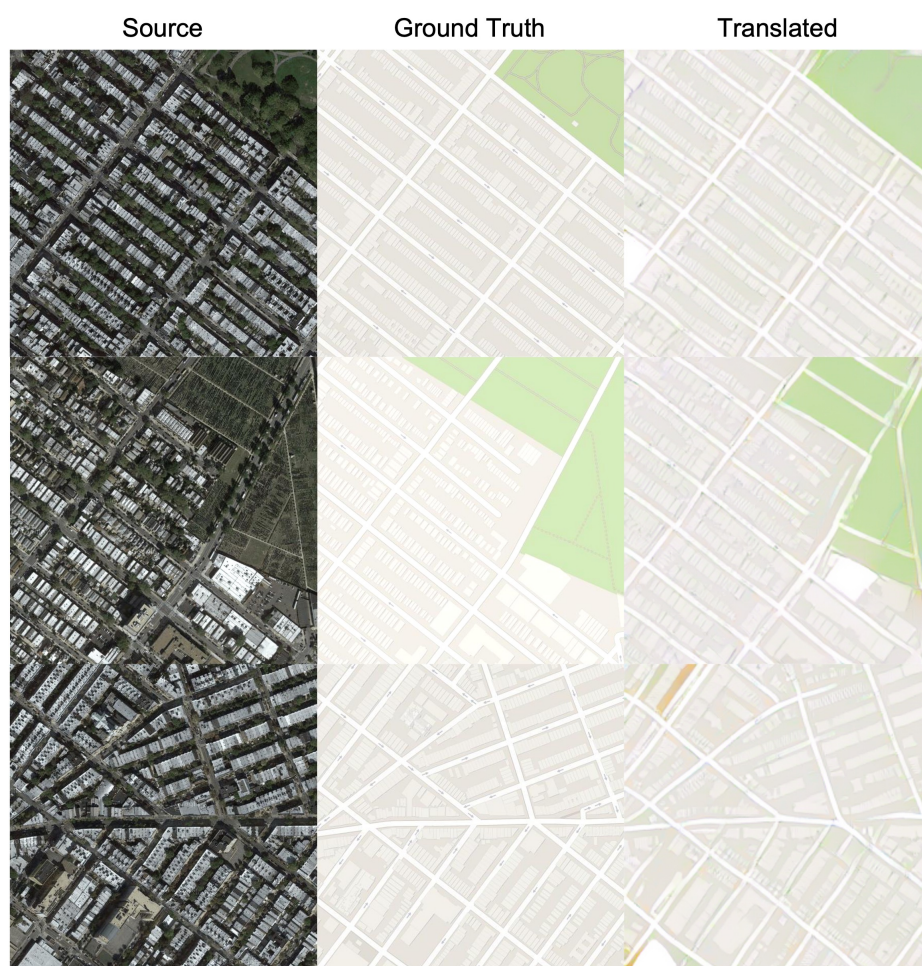Fig. 13: Example image translations using VSAIT for the Aerial Photo to Google Map experiment.

Fig. 14: Example image translations using VSAIT for the Aerial Photo to Google Map experiment.

### G.4    Visual Comparison to Baseline Methods

We additionally demonstrate the comparison between VSAIT and baseline methods for the Aerial Photo to Google Map and Google Map to Aerial Photo experiments in Figure 15. As shown in the figure, for these more difficult experiments where the datasets are specifically sub-sampled to reduce content overlap between domains, most of the baseline methods exhibit considerable semantic flipping (e.g., water to land).
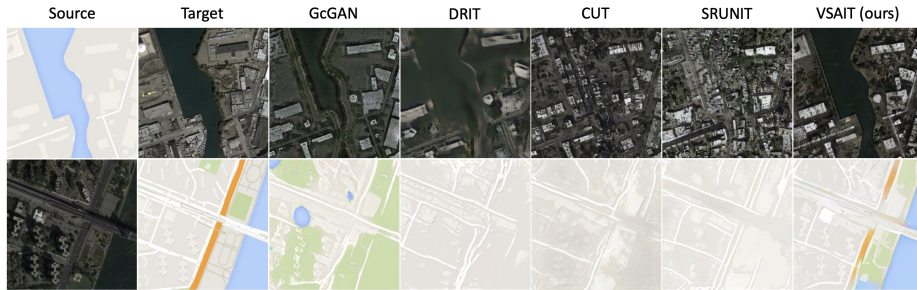


Fig. 15: Example image translations for VSAIT and baseline methods for Aerial Photo to Google Map and Google Map to Aerial Photo experiments. Note that EPE cannot be evaluated for these experiments.

### References

[1]  Jia Deng et al. "Imagenet: A large-scale hierarchical image database". In: *2009 IEEE conference on computer vision and pattern recognition*. Ieee. 2009, pp. 248–255.

[2]  Judy Hoffman et al. "Cycada: Cycle-consistent adversarial domain adaptation". In: *International conference on machine learning*. PMLR. 2018, pp. 1989–1998.

[3]  Zhiwei Jia et al. "Semantically Robust Unpaired Image Translation for Data with Unmatched Semantics Statistics". In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2021, pp. 14273–14283.

[4]  Peer Neubert, Stefan Schubert, and Peter Protzel. "An introduction to hyperdimensional computing for robotics". In: *KI-Künstliche Intelligenz* 33.4 (2019), pp. 319–330.

[5]  Karen Simonyan and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition". In: *arXiv preprint arXiv:1409.1556* (2014).

[6]  K. Wang, K. Akash, and T. Misu. "Learning Temporally and Semantically Consistent Unpaired Video-to-video Translation Through Pseudo-Supervision From Synthetic Optical Flow". In: *arXiv* (2022).

[7]   Y. Yang and S. Soatto. "FDA: Fourier domain adaptation for semantic segmentation". In: *CVPR*. 2020.

[8]   Jun-Yan Zhu et al. "Unpaired image-to-image translation using cycle-consistent adversarial networks". In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 2223–2232.