

**NANYANG
TECHNOLOGICAL
UNIVERSITY**
SINGAPORE

**AI6121 Project:
Image Translation and UDA**

**Rao Yixi, G2302775D
Jin Zhixiao, G2303771H
Huang Siwei, G2304149C**

School Of Computer Science And Engineering

18/11/2023

Contents

1	Introduction	2
2	Image-to-Image Translation	3
2.1	Introduction	3
2.2	Model Architecture	3
2.3	Implementation	4
2.3.1	Dataset & Pre-processing	4
2.3.2	Discriminator Model	5
2.3.3	Generator Model	5
2.3.4	Training	6
2.4	Result Discussion	8
2.4.1	Result presentation	8
2.4.2	Major constraints of the CycleGAN	8
3	Unsupervised Domain Adaptation via I2I Translation	12

1 Introduction

2 Image-to-Image Translation

The image to image translation model we chose is Cycle Generative Adversarial Network (CycleGAN) developed by Zhu et al [1].

2.1 Introduction

For image-to-image translation task, an original image set S and a target image set T are required. For most of the current I2I models, the two image sets used for training need to be paired, meaning that for every image $s \in S$, there must be a corresponding image $t \in T$. However, this kind of paired image set is scarce because (1) Manual production is time-consuming and laborious (2) The target image set required by some vision and graphics tasks does not exist, such as image translation between painters of different styles. The CycleGAN model proposed by Zhu et al. is designed to solve the unpaired image-to-image translation problem.

The CycleGAN model needs to learn a mapping $G : S \rightarrow T$ between an input image and an output image, which can take the special features of the original image set, and then figure out how to transform these features to the target image set, making it difficult to distinguish between the translated image distribution and the target image set distribution. However, this kind of unpaired unidirectional I2I translation will have an under-constrained problem, that is, the translated image distribution $G(s)$ does match the target image set, but it cannot correspond to the original image s in a meaningful way. This is because there are many mappings G that can generate images with the same distribution. Moreover, this unidirectional training also triggers the well-known problem of mode collapse.

Zhu et al. used the property that translation should be “cycle consistent” to define another mapping $F : T \rightarrow S$ to reproject the translated image onto the original image, ensuring that $F(G(s)) = s$ and $G(F(t)) = t$. The model uses two generators and their corresponding discriminators using adversarial loss and defines a novel cycle consistency loss to train the two generators, and discriminators simultaneously. Figure 1 shows the basic design of CycleGAN and the Loss definition.

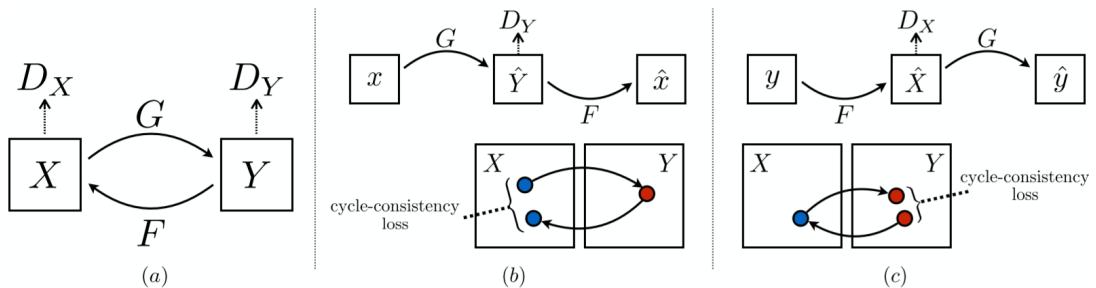


Figure 1: (a) The model contains two mapping functions (generators) and associated adversarial discriminators (b) forward cycle-consistency loss (c) backward cycle-consistency loss.

2.2 Model Architecture

To accomplish I2I, the overall model contains two generators $G : S \rightarrow T$ and $F : T \rightarrow S$ and the corresponding two adversarial discriminators D_S and D_T , which are used to distinguish

between the original image and the translated image. The basic idea is to improve the accuracy of the adversarial discriminators on image sets S or T while improving the image translation ability of the generators. The ideal result is that the discriminators maintain high accuracy on the image sets S or T , but only 50% accuracy in recognizing the fake images translated by the generators.

Adversarial Loss. Adversarial loss is applied to generators and discriminators and is mainly used to ensure that the levels of sets of translations from domain S (T) and domain T (S) are appropriate. Below is the adversarial loss for G and D_T . The other one is similar.

$$\mathcal{L}_{GAN}(G, D_T, S, T) = \mathbb{E}_{t \sim p_{data(t)}} [\log D_T(t)] + \mathbb{E}_{s \sim p_{data(s)}} [\log(1 - D_T(G(s)))] \quad (1)$$

Cycle Consistency Loss. To ensure cycle consistency, a reproject error-like forward cycle consistency loss and backward cycle consistency loss, corresponding to two generators respectively, are used, which are defined as:

$$\mathcal{L}_{cyc}(G, F) = \mathbb{E}_{s \sim p_{data(s)}} [| | | F(G(s)) - s | |_1] + \mathbb{E}_{t \sim p_{data(t)}} [| | | G(F(t)) - t | |_1] \quad (2)$$

Therefore, the full objective function can be defined as:

$$G^*, F^* = \arg \min_{G, F} \max_{D_S, D_T} \mathcal{L}(G, F, D_T, D_S) \quad (3)$$

$$\mathcal{L}(G, F, D_T, D_S) = \mathcal{L}_{GAN}(G, D_T, S, T) + \mathcal{L}_{GAN}(G, D_S, T, S) + \lambda \mathcal{L}_{cyc}(G, F) \quad (4)$$

2.3 Implementation

We have organized the implementation of CycleGAN into four steps. The first step is to define the dataset and pre-processing, the second step is to define the discriminator model, the third step is to define the generator model, and the fourth step is to define the training for CycleGAN. We use Pytorch to implement CycleGAN and refer to the source code of Zhu et al [1] and Aladdin's code reproduction [2]. Here we only present the core codes, please check the source code for detailed implementation.

2.3.1 Dataset & Pre-processing

We use the GTA5 video game street view from [3] as the original dataset S and the real cityscape dataset from [4] as the T dataset, so our goal is to convert the GTA5 video game street view image into a real street view style image. The image augmentation we use is, to resize the image to 180×360 , then use horizontal flip on the image with $p = 0.5$, and finally normalize all the image channels to have a mean of 0.5 and standard deviation of 0.5. Finally, convert it to tensor. The detailed implementation is shown below.

```
A.Compose([
    A.Resize(width=360, height=180),
    A.HorizontalFlip(p=0.5),
    A.Normalize(mean=[0.5, 0.5, 0.5],
               std=[0.5, 0.5, 0.5],
               max_pixel_value=255),
    ToTensorV2()])
])
```

2.3.2 Discriminator Model

The implementation of discriminator architectures follows the paper description of Zhu et al., who use 70×70 PatchGAN as the discriminator model. Define C_k as 4×4 Convolution-InstanceNorm-LeakyReLU layer with k filters and stride 2. We use Pytorch's `nn.Module` to define C_k block class:

```
self.conv = nn.Sequential(
    nn.Conv2d(
        in_channels,
        out_channels,
        4,
        stride,
        1,
        bias=True,
        padding_mode="reflect"
    ),
    nn.InstanceNorm2d(out_channels),
    nn.LeakyReLU(0.2, inplace=True)
)
```

We then connect four C_k blocks in the order of $C64-C128-C256-C512$ as the discriminator architecture, where the $C64$ block does not use the InstanceNorm layer. Finally we use a convolution to produce a 1-dimensional output and add the sigmoid function.

2.3.3 Generator Model

The Generator architecture used by Zhu et al. consists of a convolution block and a residual block, where the convolution block is divided into downsampling and upsampling blocks. In the following, we use `nn.Module` define convolution block as:

```
# Downsampling layers or upsampling layers
self.conv = nn.Sequential(
    nn.Conv2d(in_channels,
              out_channels,
              padding_mode="reflect",
              **kwargs)
    if down
    else
        nn.ConvTranspose2d(in_channels,
                          out_channels,
                          **kwargs),
    nn.InstanceNorm2d(out_channels),
    nn.ReLU(inplace=True) if use_act else nn.Identity()
)
```

Based on this convolution block class, we can define three types of convolution blocks. The first is the $c7s1-k$ block, which uses a 7×7 Convolution-InstanceNorm-ReLU layer with k filters and stride 1. The second is the downsampling block dk , which uses a 3×3 Convolution-InstanceNorm-ReLU layer with k filters and stride 2. The third is the upsampling block uk , which is a 3×3 fractional-strided-Convolution-InstanceNorm-ReLU layer with k filters and stride $\frac{1}{2}$.

For residual block, it is two 3×3 convolutional layers with the same number of filters on both layer. We defined it as the residual block (below) class using `nn.Module`. The `forward` function is `return x + self.block(x)`.

```
self.block = nn.Sequential(
    ConvBlock(channels,
              channels,
              kernel_size=3,
              padding=1),
    ConvBlock(channels,
              channels,
              use_act=False,
              kernel_size=3,
              padding=1)
)
```

We build the generator network as: “c7s1-64,d128,d256,R256,R256,R256,R256, R256,R256, R256, R256, R256, R256, R256, R256, R256, R256, R256, u128 u64,c7s1-3”, and followed by a tanh activation function.

```
# c7s1-64
x = self.initial(x)
# d128, d256
for layer in self.down_blocks:
    x = layer(x)
# R256, R256, R256, R256, R256, R256, R256, R256, R256, R256
x = self.res_blocks(x)
# u128 u64
for layer in self.up_blocks:
    x = layer(x)
# c7s1-3
return torch.tanh(self.last(x))
```

2.3.4 Training

When implementing training, we adopted Zhu et al.’s training details and their recommended settings, such as training parameter settings:

- Use a batch size of 2 to prevent out of memory.
- A total of 200 epochs are trained, of which the learning rate is kept constant for 100 epochs, and the decay learning rate strategy is used for the other 100 epochs. Here we use the cosine annealing learning rate strategy.
- Set the initial learning rate to 0.0002
- Set cycle consistent control factor $\lambda = 10$
- Use Adam solver

At the same time, we adopt their recommended improvement techniques, such as using least-squares loss to replace the negative log likelihood objective in \mathcal{L}_{GAN} . This loss is more stable during training and generates higher quality results. Therefore, the new goal will be:

- When training G or F : minimize

$$\mathbb{E}_{s \sim p_{data(s)}}[(D_T(G(s)) - 1)^2] +$$

$$\mathbb{E}_{t \sim p_{data(t)}}[(D_S(F(t)) - 1)^2] +$$

$$\mathcal{L}_{cyc}(G, F)$$

- When training D_S or D_T : minimise

$$\mathbb{E}_{s \sim p_{data(s)}}[(D_S(s) - 1)^2] + \mathbb{E}_{t \sim p_{data(t)}}[D_S(F(t))^2] +$$

$$\mathbb{E}_{t \sim p_{data(t)}}[(D_T(t) - 1)^2] + \mathbb{E}_{s \sim p_{data(s)}}[D_T(G(s))^2]$$

Following the above training setting, for each batch images (source, target), we first train the discriminators D_T and D_S . We can obtain two discriminators' result by using corresponding generators.

```
# Patch: DT(T) and DT(GT(S))
DT_real = disc_T(target)
DT_fake = disc_T(gen_T(source))

# Patch: DS(S) and DS(GS(T))
DS_real = disc_S(source)
DS_fake = disc_S(gen_S(target))
```

And then we can calculate the \mathcal{L}_{GAN} losses with respect to discriminators D_T and D_S .

```
# L_GAN(GT, DT, S, T) = minimize Et[(DT(t) - 1)^2] + Es[DT(GT(s))^2]
DT_real_loss = mse(DT_real, torch.ones_like(DT_real))
DT_fake_loss = mse(DT_fake, torch.zeros_like(DT_fake))
# L_GAN(GT, DT, S, T)
DT_loss = DT_real_loss + DT_fake_loss

# L_GAN(GS, DS, S, T) = minimize Es[(DS(s) - 1)^2] + Et[DS(GS(t))^2]
DS_real_loss = mse(DS_real, torch.ones_like(DS_real))
DS_fake_loss = mse(DS_fake, torch.zeros_like(DS_fake))
# L_GAN(GS, DS, S, T)
DS_loss = DS_real_loss + DS_fake_loss
```

Finally, we put it together to get the total adversarial loss, and then train the model with the defined optimizer to optimize discriminators.

```
# put it together: L_GAN(GT, DT, S, T) + L_GAN(GS, DS, S, T)
D_loss = (DT_loss + DS_loss) / 2
```

Next, we need to train the generators G_t and G_s . In the full objective function, the two generator models appear both in adversarial loss and cycle consistency loss. To compute the adversarial loss in terms of generators, the results of the discriminators are first computed for the forward and backward translated images.

```
# Patch: DT(GT(S)) and DS(GS(T))
DT_fake = disc_T(gen_T(source))
DS_fake = disc_S(gen_S(target))
```

Then we need to compute the adversarial loss with respect to generators.

```
# minimize Et [(DT(t) - 1)^2] and Es [(DS(s) - 1)^2]
loss_GT = mse(DT_fake, torch.ones_like(DT_fake))
loss_GS = mse(DS_fake, torch.ones_like(DS_fake))
```

To compute the cycle consistency loss, which only concludes the generators, we first need to get the reprojeciton of the translated images.

```
# GS(GT(S)) -> S'
cycle_source = gen_S(gen_T(source))
# GT(GS(T)) -> T'
cycle_target = gen_T(gen_S(target))
```

And then calculate the cycle consistency loss according to the definition shown above. Finally, minimising the cycle consistency loss by using the defined optimizer to optimize the generators model.

```
# Lcyc(GS,GT) = Es [|GS(GT(s)) - s|] + Et [|GT(GS(t)) - t|]
cycle_source_loss = l1(source, cycle_source)
cycle_target_loss = l1(target, cycle_target)

# add all together: L(GS,GT,DS,DT) = LGAN(GT,DT,X,Y) + LGAN(GS,DS,X,Y) +
#                                         Lcyc(GS,GT)
G_loss = loss_GT + loss_GS +
lambda * (cycle_source_loss + cycle_target_loss)
```

2.4 Result Discussion

2.4.1 Result presentation

We randomly sampled 3000 images from GTA5 image set and 3000 images from cityscape image set for training and continued to randomly sample 1000 images from the remaining GTA5 image set for testing. Figure 2 are the results of the test which contains the testing results from the training set and the results from the test set. We can observe that our CycleGAN recognizes color and texture variations excellently across different image sets and successfully converts video game cityscapes to real-world style cityscapes. It performs better on the training set than on the test set, which conforms to our expectations. Note that the converted image on the test set seems to be slightly blurry compared to the image on the training set, and there may be some color distortion.

Figure 3 shows CycleGAN’s byproduct generator F , which converts real-world cityscape into video game style.

2.4.2 Major constraints of the CycleGAN

Wrong recognition of objects. In Figure 4, we find that CycleGAN has a constraint that recognises the wrong class of objects. For example, it may convert the clouds in the source image into trees in the translated image, or in Figure 4 a, it converts the clouds in the source image into buildings in the translated image, and it may also interpret the cloud above the tree in the source image as a portion of the tree, and then generate a huge tree, such as in Figure 4 c. The reason for this phenomenon may be that CycleGAN is better at achieving color



Figure 2: Our CycleGAN model evaluation results on training set (left) and testing set (right).



Figure 3: The CycleGAN model’s byproduct: backward translation results.

and texture changes rather than complex geometric changes, especially clouds and trees with irregular shapes. This limitation was also mentioned by Zhu et al. in [1].

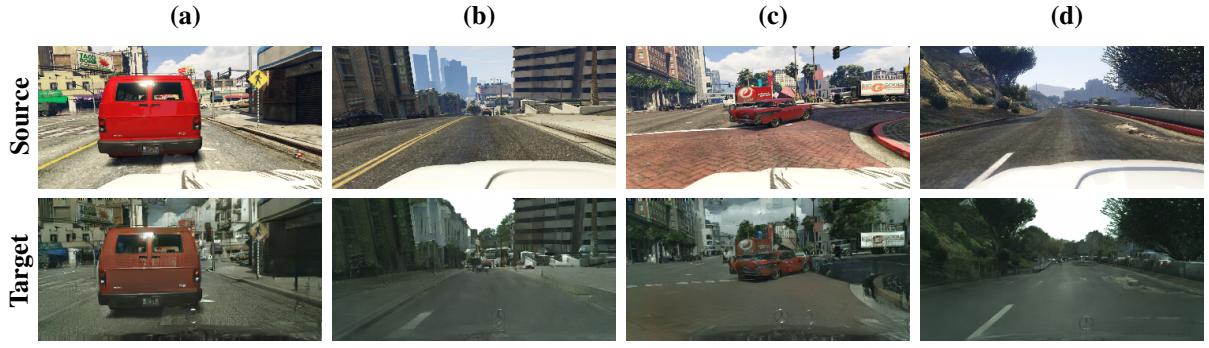


Figure 4: Failure cases of the CycleGAN: Wrong recognition of the class of the object.

Generation of artifact. In Figure 5 we find that CycleGAN may generate some artifacts that do not exist in the real world, such as in the translated images in examples a, b, and c in Figure 5. A strange dark green patch is generated in the sky. Besides, in example b in Figure 5, in the right part of the image, the image which was originally a sunny sky was translated into an image with a dark cloud artifact. I think the reason for this phenomenon may be the poor performance of the learned discriminators, which may mislead the generator, and thus the image with artifacts may be mistakenly recognized as the correct one when performing reprojection during the training process. As a result, balancing the training speed of discriminators and generators and balancing the adversarial loss and the cycle consistency loss is also a constraint for CycleGAN.

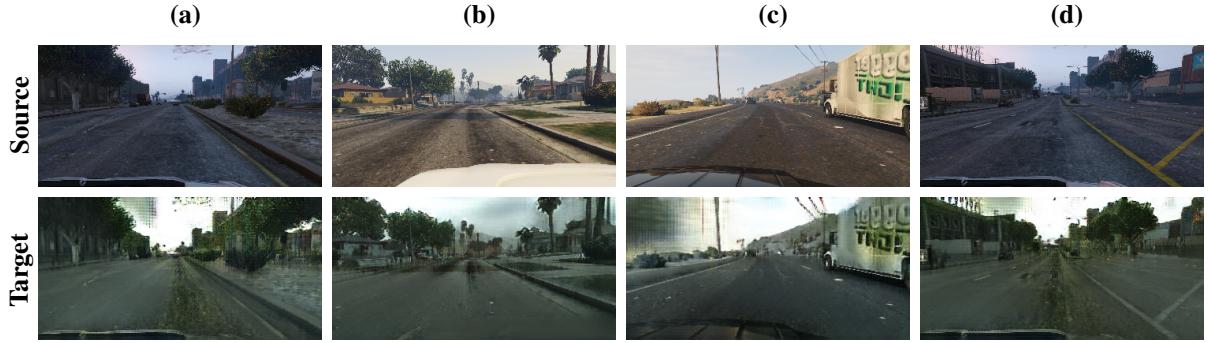


Figure 5: Failure cases of the CycleGAN: Wrong generation of artifact.

Generation of artifact. In Figure 6, we observe that CycleGAN performs poorly in night and dusk scenes. In the night scene, it translates the dark sky into the hue that does not exist in real-world cityscapes. In the dusk scene, it translates the orange sky into a sunny sky, failing to accurately capture the temporal information of the original image. The reason for this phenomenon may be the fact that the training set does not contain such images, or there are a few pairs of such training images compared to daytime images. The generalizability of CycleGAN is also very weak compared to other models using the aligned paired training method.

Perform badly at night and dusk scenes. In Figure 6, we observe that CycleGAN performs poorly in night and dusk scenes. In the night scene, it translates the dark sky into the hue that does not exist in real-world cityscapes. In the dusk scene, it translates the orange sky

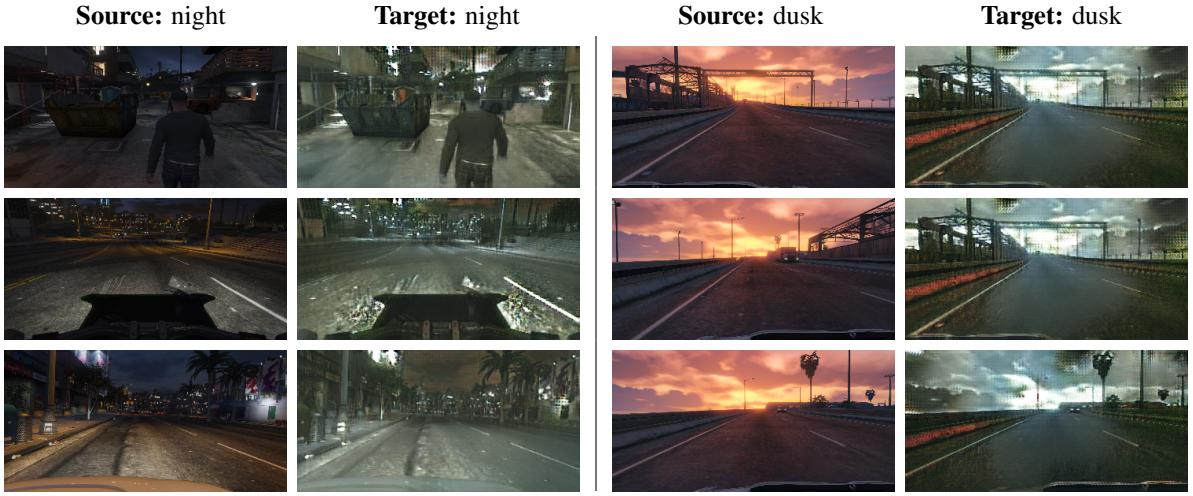


Figure 6: Failure cases of the CycleGAN: Performs poorly in night and dusk scenes.

into a sunny sky, failing to accurately capture the temporal information of the original image. The reason for this phenomenon may be the fact that the training set does not contain such images, or there are a few pairs of such training images compared to daytime images. The generalizability of CycleGAN is also very weak compared to other models using the aligned paired training method.

The fuzzy images and Color distortion. In the above Figures 4, 5, and 6, we can see that the images generated on the test set have less clarity compared to the source images and there will be some color distortion. The reason for this phenomenon is that we used too few images in the training set or too few epochs. Thus we can conclude that although CycleGAN can learn mappings without paired images, it still needs a lot of data to achieve good results in some cases, especially when the target area changes greatly.

Intrinsic constraints of CycleGAN. The mapping generated by CycleGAN is a kind of deterministic mapping because cycle consistency enforces G , and F to be inverses of each other [5]. When dealing with complex cross-domain datasets, this deterministic mapping causes CycleGAN to learn an arbitrary one-to-one mapping instead of capturing the true, structured conditional distribution more faithfully [5]. At the same time, CycleGAN still does not eliminate the mode collapse problem. It still exists that the generator tends to generate similar or identical outputs while ignoring some subtle differences in the inputs. This may lead to a lack of diversity in the generated results.

3 Unsupervised Domain Adaptation via I2I Translation

References

- [1] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, “Unpaired image-to-image translation using cycle-consistent adversarial networks,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [2] A. Persson, “A clean, simple and readable implementation of cyclegan in pytorch,” <https://github.com/aladdinpersson/Machine-Learning-Collection/tree/master/ML/Pytorch/GANs/CycleGAN>, accessed: 2023-11-10.
- [3] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, “Playing for data: Ground truth from computer games,” 2016.
- [4] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, “The cityscapes dataset for semantic urban scene understanding,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3213–3223.
- [5] A. Almahairi, S. Rajeshwar, A. Sordoni, P. Bachman, and A. Courville, “Augmented cycle-gan: Learning many-to-many mappings from unpaired data,” in *International conference on machine learning*. PMLR, 2018, pp. 195–204.