



UNIVERSITÉ PARIS SACLAY

ADVANCED IN MACHINE VISION

MASTER OF RESEARCH IN ARTIFICIAL INTELLIGENCE AND
MACHINE VISION

Lab2

Image-to-Image Translation for Semantic Segmentation

Student:
Nguyen Trong Dat

Supervisor:
Hedi Tabia

Master:
M2MMVAI

Academic Year
2023-2024

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 5 |
| 1.1 | Cityscape Dataset: A Brief Overview | 5 |
| 1.1.1 | Loading the Dataset | 5 |
| 1.1.2 | Visualization of Sample Images | 5 |
| 1.2 | Objective of the Study | 5 |
| 2 | CGAN: Our model | 6 |
| 2.1 | GAN | 6 |
| 2.2 | CGAN: Specifics | 6 |
| 2.3 | Implementation of Conditional GAN (CGAN) | 6 |
| 2.3.1 | Generator Implementation | 6 |
| 2.3.2 | Discriminator Implementation | 7 |
| 2.4 | Training of Conditional GAN (CGAN) | 8 |
| 2.5 | Results | 9 |
| 2.5.1 | Pixel-wise Accuracy | 9 |
| 2.5.2 | Results | 9 |
| 2.5.3 | Sample Result Image | 10 |
| 3 | U-Net : Baseline-model | 10 |
| 3.1 | Introduction to U-Net | 10 |
| 3.1.1 | Implementation of U-Net | 10 |
| 3.2 | Training U-Net | 12 |
| 3.3 | Results | 13 |
| 3.3.1 | Pixel-wise Accuracy | 13 |
| 3.3.2 | Sample Result Image | 13 |
| 4 | Comparison: U-Net vs. CGAN | 13 |
| 4.1 | Accuracy Comparison | 13 |
| 4.2 | Generated Images Comparison | 14 |
| 4.3 | Comment | 14 |

List of Figures

| | | |
|---|---|----|
| 1 | Cityscape dataset | 5 |
| 2 | Comparison of Generated Image and Ground Truth | 10 |
| 3 | Comparison of Generated Image and Ground Truth (Unet) | 13 |
| 4 | Comparison of Generated Images - U-Net (top) vs CGAN (bottom) | 14 |

Listings

| | | |
|---|---|----|
| 1 | Python code for Loading Dataset | 5 |
| 2 | Python code for Generator | 7 |
| 3 | Python code for Discriminator | 7 |
| 4 | Python code for Training the Model | 8 |
| 5 | Python code for Convolutional Block | 10 |
| 6 | Python code for Encoder and Decoder | 11 |
| 7 | Python code for Unet | 11 |
| 8 | Python code for Unet | 12 |

List of Tables

| | | |
|---|--|----|
| 1 | Pixel-wise Accuracy on Training and Validation Sets | 9 |
| 2 | Pixel-wise Accuracy on Training and Validation Sets | 13 |
| 3 | Comparison of Pixel-wise Accuracy between U-Net and CGAN | 13 |

1 Introduction

In this report, we delve into the realm of semantic segmentation using the Cityscape dataset, a benchmark dataset designed for urban scene understanding. The Cityscape dataset provides a diverse collection of high-quality images captured across various cities, making it an invaluable resource for training and evaluating computer vision models, especially in the context of autonomous driving and urban environment analysis.

1.1 Cityscape Dataset: A Brief Overview

The Cityscape dataset comprises images taken from urban street scenes, annotated with pixel-level semantic labels. These annotations include detailed information about objects and regions within the images, such as roadways, vehicles, pedestrians, and more. This dataset serves as a pivotal benchmark for advancing research in semantic segmentation and related computer vision tasks.

1.1.1 Loading the Dataset

```

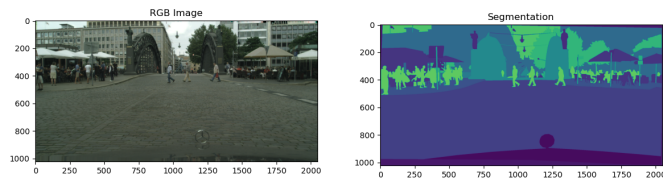
1 class CityscapeDataset(Dataset):
2     def __init__(self, image_path, transform_img=None, transform_label=None):
3         self.image_path = image_path
4         self.transform_img = transform_img
5         self.transform_label = transform_label
6
7     def __len__(self):
8         return len(self.image_path)
9
10    def __getitem__(self, idx):
11        img = plt.imread(self.image_path[idx])
12
13        image, label = img[:, :img.shape[1] // 2], img[:, img.shape[1] // 2:]
14
15        if self.transform_img:
16            image = self.transform_img(image)
17
18        if self.transform_label:
19            label = self.transform_label(label)
20
21        return image, label

```

Listing 1: Python code for Loading Dataset

1.1.2 Visualization of Sample Images

Let's visualize one pair of original and ground truth images from the Cityscape dataset to better understand the data. The following figure shows an example:



(a) Original Image

(b) Ground truth Image

Figure 1: Cityscape dataset

1.2 Objective of the Study

Our primary objective is to explore and compare the performance of two distinct approaches to semantic segmentation: a conventional model, exemplified by architectures like VAE (Variational

Autoencoder) and U-Net, and a Conditional Generative Adversarial Network (CGAN). Through this comparative analysis, we aim to gain insights into the strengths and limitations of each approach in the specific context of urban scene segmentation.

By leveraging the Cityscape dataset, we aim to assess the models' ability to accurately label and segment different objects and entities within complex urban scenes. The utilization of both traditional and generative adversarial approaches allows us to evaluate the trade-offs between model complexity, computational efficiency, and the quality of segmentation outputs.

This study not only contributes to the ongoing discourse in the field of computer vision but also provides practical insights into the applicability of different models for semantic segmentation in urban environments. The following sections will detail the methodology, experimental setup, results, and conclusions derived from our exploration of these models on the Cityscape dataset.

2 CGAN: Our model

2.1 GAN

A Generative Adversarial Network (GAN) is a class of machine learning models introduced by Goodfellow et al. in 2014. The GAN framework consists of two main components: a generator (G) and a discriminator (D). The primary objective of a GAN is to generate synthetic data that is indistinguishable from real data.

- **Generator (G):** Takes random noise as input and generates synthetic data.
- **Discriminator (D):** Takes both real and generated data as input and attempts to distinguish between them.

During training, G aims to create realistic data to fool D , while D improves its ability to differentiate between real and generated samples. This adversarial process results in the generator producing increasingly convincing data over time.

2.2 CGAN: Specifics

A Conditional Generative Adversarial Network (CGAN) extends the GAN framework by introducing conditional information to guide the data generation process. In a CGAN, both G and D receive additional conditional information along with the random noise.

- **Generator (G):** Takes random noise and conditional information to generate data with specific characteristics.
- **Discriminator (D):** Receives both real and generated data along with conditional information to make more informed decisions.

The conditional information allows for targeted data generation, making CGANs suitable for tasks where specific attributes or classes are desired in the generated samples.

2.3 Implementation of Conditional GAN (CGAN)

To further explore the application of Conditional Generative Adversarial Networks (CGAN) in the context of semantic segmentation, we present the implementation of a CGAN tailored for our task. The CGAN consists of a generator (G) and a discriminator (D), each designed to accommodate conditional information for targeted image generation.

2.3.1 Generator Implementation

The generator is responsible for synthesizing realistic images based on random noise and conditional information. Below is the Python code for the CGAN generator:

```

1 class Generator(nn.Module):
2     def __init__(self, ngf, img_size, nc):
3         super().__init__()
4         self.nc = nc
5         self.img_size = img_size
6         self.ngf = ngf
7
8         self.downSampling = nn.Sequential(
9             nn.Conv2d(nc, self.ngf, 4, 2, 1, bias=False),
10            nn.LeakyReLU(0.2, inplace=True),
11
12            nn.Conv2d(self.ngf, self.ngf * 2, 4, 2, 1, bias=False),
13            nn.BatchNorm2d(self.ngf * 2),
14            nn.LeakyReLU(0.2, inplace=True),
15
16            nn.Conv2d(self.ngf * 2, self.ngf * 4, 4, 2, 1, bias=False),
17            nn.BatchNorm2d(self.ngf * 4),
18            nn.LeakyReLU(0.2, inplace=True),
19
20            nn.Conv2d(self.ngf * 4, self.ngf * 8, 4, 2, 1, bias=False),
21            nn.BatchNorm2d(self.ngf * 8),
22            nn.LeakyReLU(0.2, inplace=True),
23
24            nn.Conv2d(self.ngf * 8, self.ngf * 16, 4, 1, 0, bias=False),
25        )
26
27        self.upSampling = nn.Sequential(
28            nn.ConvTranspose2d(self.ngf * 16, self.ngf * 16, 4, 1, 0, bias=False),
29            nn.BatchNorm2d(self.ngf * 16),
30            nn.ReLU(True),
31
32            nn.ConvTranspose2d(self.ngf * 16, self.ngf * 8, 4, 2, 1, bias=False),
33            nn.BatchNorm2d(self.ngf * 8),
34            nn.ReLU(True),
35
36            nn.ConvTranspose2d(self.ngf * 8, self.ngf * 4, 4, 2, 1, bias=False),
37            nn.BatchNorm2d(self.ngf * 4),
38            nn.ReLU(True),
39
40            nn.ConvTranspose2d(self.ngf * 4, self.ngf * 2, 4, 2, 1, bias=False),
41            nn.BatchNorm2d(self.ngf * 2),
42            nn.ReLU(True),
43
44            nn.ConvTranspose2d(self.ngf * 2, nc, 4, 2, 1, bias=False),
45            nn.Tanh()
46        )
47
48        def forward(self, input):
49            x = self.downSampling(input)
50            out = self.upSampling(x)
51            return out

```

Listing 2: Python code for Generator

The generator architecture is structured to take into account both the random noise and conditional information, producing images that align with specific characteristics.

2.3.2 Discriminator Implementation

The discriminator evaluates the authenticity of input images, considering both real and generated samples, along with the associated conditional information. Here is the Python code for the CGAN discriminator:

```

1 class Discriminator(nn.Module):
2     def __init__(self, ngpu, img_size, nc):
3         super(Discriminator, self).__init__()
4         self.ngpu = ngpu
5         self.img_size = img_size
6         self.nc = nc
7

```



```

8     self.main = nn.Sequential(
9
10         nn.Conv2d(nc*2, ndf, 4, 2, 1, bias=False),
11         nn.LeakyReLU(0.2, inplace=True),
12
13         nn.Conv2d(ndf, ndf * 2, 4, 2, 1, bias=False),
14         nn.BatchNorm2d(ndf * 2),
15         nn.LeakyReLU(0.2, inplace=True),
16
17         nn.Conv2d(ndf * 2, ndf * 4, 4, 2, 1, bias=False),
18         nn.BatchNorm2d(ndf * 4),
19         nn.LeakyReLU(0.2, inplace=True),
20
21         nn.Conv2d(ndf * 4, ndf * 8, 4, 2, 1, bias=False),
22         nn.BatchNorm2d(ndf * 8),
23         nn.LeakyReLU(0.2, inplace=True),
24
25         nn.Conv2d(ndf * 8, 1, 4, 1, 0, bias=False),
26         nn.Sigmoid()
27     )
28
29
30     def forward(self, x, labels):
31         x = torch.cat([x, labels], 1)
32         out = self.main(x)
33         return out.view(-1,1,1,1)

```

Listing 3: Python code for Discriminator

The discriminator is designed to make more informed decisions by considering both image content and conditional information.

This CGAN implementation will be utilized in our experiments to assess its effectiveness in semantic segmentation tasks using the Cityscape dataset. In the subsequent sections, we will detail the experimental setup, results, and conclusions derived from applying the CGAN model to our specific use case.

2.4 Training of Conditional GAN (CGAN)

The training loop for the Conditional Generative Adversarial Network (CGAN) involves an adversarial process between the generator and discriminator. The goal is to train the generator to produce realistic images that the discriminator cannot distinguish from real ones. The training code is presented below:

```

1 print("Starting Training Loop...")
2
3 n_of_epochs = 200
4 loss_Discriminator = 0
5
6 for e in range(n_of_epochs):
7     running_loss_D = 0.0
8     print('Starting epoch {}'.format(e+1))
9     for i, (input,output) in enumerate(train_loader, 0):
10
11         #####
12         ## Update Discriminator: maximize log(D(x/y)) + log(1 - D(G(z/y))) ##
13         #####
14
15         # Loss on real Images
16         myDiscriminator.zero_grad()
17         predicted = myDiscriminator(input.to(device),output.to(device))
18         loss_real = criterion(predicted, torch.ones((len(predicted)), 1, 1,1, dtype
19 =torch.float, device=device).to(device))
20         loss_real.backward(retain_graph=True)
21
22         # Loss on fake Images
23         fake_output = myGenerator(input.to(device))
24         predicted = myDiscriminator(input.to(device),fake_output.to(device))

```

```

24     loss_fake = criterion(predicted, torch.zeros((len(predicted)), 1, 1, 1,
dtype=torch.float, device=device).to(device))
25     loss_fake.backward(retain_graph=True)
26     loss_Discriminator = loss_real + loss_fake
27     d_optimizer.step()
28
29     #####
30     ## Update Generator: maximize log(D(G(z/y))) ##
31     #####
32
33     myGenerator.zero_grad()
34     predicted = myDiscriminator(input.to(device), fake_output.to(device))
35     loss_gen = criterion(predicted, torch.ones((len(predicted)), 1, 1, 1, dtype=
torch.float, device=device).to(device))
36     loss_gen.backward()
37     g_optimizer.step()
38     running_loss_D += loss_Discriminator.item()
39
40     print('g_loss: {}, d_loss: {}'.format(loss_gen, loss_Discriminator))

```

Listing 4: Python code for Training the Model

The training loop consists of two main steps: updating the discriminator to distinguish between real and fake images, and updating the generator to produce more convincing images. These steps are repeated for a specified number of epochs. The loss values for the generator and discriminator are printed for each epoch to monitor the training process.

This training procedure aims to find a balance between the generator and discriminator, resulting in a generator capable of producing synthetic images that align with the conditional information provided.

2.5 Results

As a metric for evaluating the performance of our Conditional Generative Adversarial Network (CGAN) in semantic segmentation, we employ pixel-wise accuracy. Pixel-wise accuracy measures the proportion of correctly classified pixels compared to the total number of pixels in the images.

2.5.1 Pixel-wise Accuracy

Pixel-wise accuracy is calculated using the following formula:

$$PixelwiseAccuracy = \frac{NumberOfCorrectlyClassifiedPixels}{TotalNumberOfPixels} \times 100 \quad (1)$$

This metric provides insights into how well the model is able to correctly classify each pixel in the generated images.

2.5.2 Results

| Dataset | Pixel-wise Accuracy (%) |
|----------------|-------------------------|
| Training Set | 46.41 |
| Validation Set | 37.26 |

Table 1: Pixel-wise Accuracy on Training and Validation Sets

Although these percentages may seem relatively low, visually inspecting the segmented images reveals promising results. The apparent discrepancy between accuracy percentages and visual quality can be attributed to the presence of noise in the generated images.

It's important to note that the model's performance could further improve with additional training epochs. The current results, while demonstrating the potential of the CGAN for semantic segmentation, indicate that more training may lead to a reduction in noise and an enhancement of overall segmentation quality.

2.5.3 Sample Result Image

Let's visualize a sample result image from the validation set:

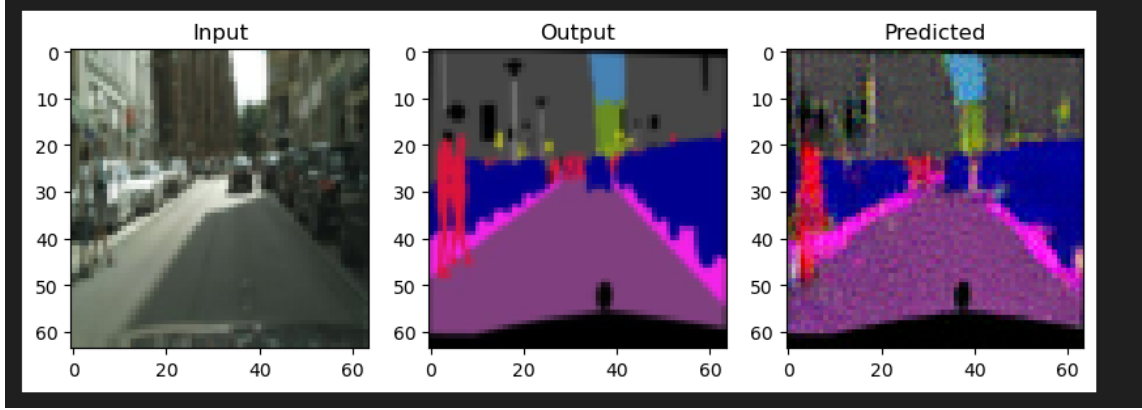


Figure 2: Comparison of Generated Image and Ground Truth

In this example, we can observe the generated image alongside the corresponding ground truth. Despite the pixel-wise accuracy, the visual comparison showcases the model's capability to capture semantic information, providing a foundation for further refinement and improvement.

3 U-Net : Baseline-model

3.1 Introduction to U-Net

U-Net is a convolutional neural network architecture designed for semantic segmentation tasks, introduced by Olaf Ronneberger, Philipp Fischer, and Thomas Brox in 2015. It is particularly well-suited for tasks where precise localization of objects in images is essential, such as medical image segmentation.

The U-Net architecture is characterized by its distinctive U-shaped structure, consisting of an encoder path and a decoder path. The encoder path captures hierarchical features from the input image through a series of convolutional and pooling layers, gradually reducing spatial resolution. The decoder path, on the other hand, upsamples the feature maps and combines them with the corresponding high-resolution features from the encoder, enabling precise localization.

One of the key innovations of U-Net is the introduction of skip connections between the encoder and decoder, allowing the model to retain detailed information during upsampling. This enables U-Net to generate accurate segmentation masks, especially in areas with fine details.

In the context of our study, U-Net serves as the baseline model for semantic segmentation on the Cityscape dataset. Its architecture and proven performance make it a valuable benchmark against which we can compare the results obtained from the Conditional Generative Adversarial Network (CGAN).

3.1.1 Implementation of U-Net

To implement U Net, we define a fundamental building block known as a convolutional block (convblock). This block consists of two convolutional layers with batch normalization and ReLU activation functions. Below is the Python code for the convblock:

```
1 class conv_block(nn.Module):
2     def __init__(self, in_c, out_c):
3         super().__init__()
4         self.conv1 = nn.Conv2d(in_c, out_c, kernel_size=3, padding=1)
5         self.bn1 = nn.BatchNorm2d(out_c)
6         self.conv2 = nn.Conv2d(out_c, out_c, kernel_size=3, padding=1)
7         self.bn2 = nn.BatchNorm2d(out_c)
8         self.relu = nn.ReLU(inplace=True)
```

```

9
10 def forward(self, x):
11     x = self.conv1(x)
12     x = self.bn1(x)
13     x = self.relu(x)
14     x = self.conv2(x)
15     x = self.bn2(x)
16     x = self.relu(x)
17     return x

```

Listing 5: Python code for Convolutional Block

The convblock is a crucial component in constructing the encoder and decoder paths of the U-Net architecture. It helps capture and process features at different scales, contributing to the network's ability to perform accurate semantic segmentation.

```

1 class encoder_block(nn.Module):
2     def __init__(self, in_c, out_c):
3         super().__init__()
4         self.conv = conv_block(in_c, out_c)
5         self.pool = nn.MaxPool2d(2, 2)
6
7     def forward(self, x):
8         x = self.conv(x) # for skip connection feature map to decoder
9         p = self.pool(x)
10        return x, p
11
12 class decoder_block(nn.Module):
13     def __init__(self, in_c, out_c):
14         super().__init__()
15         self.convT = nn.ConvTranspose2d(in_c, out_c, kernel_size=2, padding=0,
16         stride=2)
17         self.conv = conv_block(2 * out_c, out_c)
18
19     def forward(self, inputs, skips):
20         x = self.convT(inputs)
21         x = torch.cat([x, skips], axis=1)
22         x = self.conv(x)
23         return x

```

Listing 6: Python code for Encoder and Decoder

The encoderblock handles the downsampling process, consisting of a convolutional block followed by max-pooling. The decoderblock, on the other hand, performs upsampling using transpose convolution and concatenates the features with skip connections from the corresponding encoder block.

These blocks contribute to the overall architecture of U-Net, facilitating the extraction and integration of features at different scales.

To complete the implementation of U-Net, we define the final architecture by combining the previously introduced building blocks. The 'Unet' class serves as the overarching model, consisting of encoder, bottleneck, decoder, and classifier components.

```

1 class Unet(nn.Module):
2     def __init__(self):
3         super().__init__()
4         """Encoder part"""
5         self.en1 = encoder_block(3, 64)
6         self.en2 = encoder_block(64, 128)
7         self.en3 = encoder_block(128, 256)
8         self.en4 = encoder_block(256, 512)
9
10        """Bottleneck"""
11        self.b = conv_block(512, 1024)
12
13        """Decoder part"""
14        self.d1 = decoder_block(1024, 512)
15        self.d2 = decoder_block(512, 256)
16        self.d3 = decoder_block(256, 128)
17        self.d4 = decoder_block(128, 64)

```

```

18
19     """Classifier"""
20     self.outputs = nn.Conv2d(64, 3, kernel_size=1, padding=0)
21
22     def forward(self, input):
23         s1, p1 = self.en1(input)
24         s2, p2 = self.en2(p1)
25         s3, p3 = self.en3(p2)
26         s4, p4 = self.en4(p3)
27
28         b = self.b(p4)
29
30         d1 = self.d1(b, s4)
31         d2 = self.d2(d1, s3)
32         d3 = self.d3(d2, s2)
33         d4 = self.d4(d3, s1)
34
35         out = self.outputs(d4)
36
37     return out

```

Listing 7: Python code for Unet

The ‘Unet’ class orchestrates the entire U-Net architecture, consisting of encoder blocks, a bottleneck, decoder blocks, and a final classifier. This model is designed for semantic segmentation tasks and will be utilized as a baseline for comparison in our study.

3.2 Training U-Net

To train the U-Net model, we follow a standard procedure involving iterative epochs, forward and backward passes, and optimization. The code snippet below exemplifies the training process, where training and validation losses are computed and visualized at regular intervals.

```

1 train_loss = []
2 val_loss = []
3
4 for epoch in range(epochs):
5     trainloss = 0
6     valloss = 0
7     c = 0
8
9     for img, lab in tqdm(train_loader):
10         optimizer.zero_grad()
11         img = img.to(device)
12         lab = lab.to(device)
13         output = model(img)
14         loss = loss_func(output, lab)
15         loss.backward()
16         optimizer.step()
17         trainloss += loss.item()
18
19         if epoch % 5 == 0 and c < 2:
20             show([img[0], lab[0], output[0]])
21             plt.show()
22             c += 1
23
24     train_loss.append(trainloss / len(train_loader))
25
26     for img, lab in tqdm(valid_loader):
27         img = img.to(device)
28         lab = lab.to(device)
29         output = model(img)
30         loss = loss_func(output, lab)
31         valloss += loss.item()
32
33     val_loss.append(valloss / len(valid_loader))
34

```

```
35 print("epoch: {}, train loss: {}, valid loss: {}".format(epoch, train_loss[-1],  
val_loss[-1]))
```

Listing 8: Python code for Unet

This code utilizes training and validation loaders, computes the loss during training, performs backpropagation, and updates the model's parameters. Additionally, it includes a visualization of sample images at regular intervals.

3.3 Results

3.3.1 Pixel-wise Accuracy

| Dataset | Pixel-wise Accuracy (%) |
|----------------|-------------------------|
| Training Set | 56.75 |
| Validation Set | 59.88 |

Table 2: Pixel-wise Accuracy on Training and Validation Sets

3.3.2 Sample Result Image

Let's visualize a sample result image from the validation set:

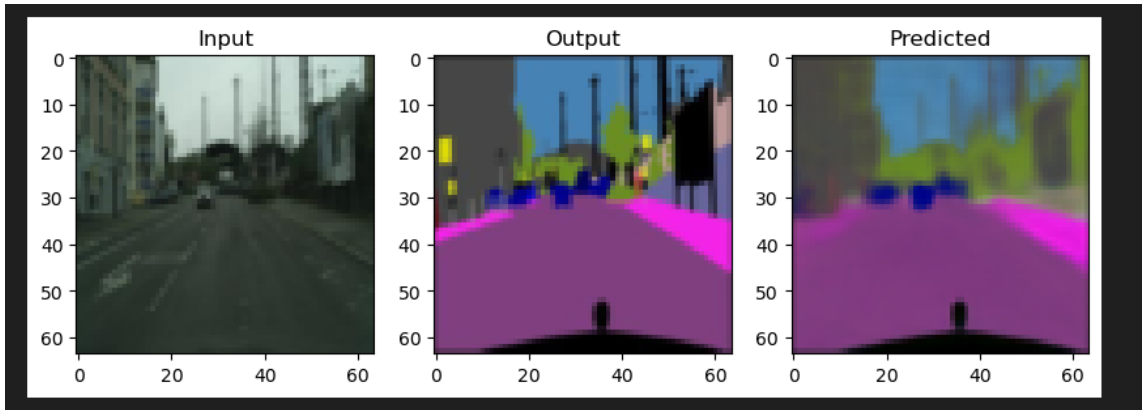


Figure 3: Comparison of Generated Image and Ground Truth (Unet)

In this example, we can observe the generated image using U-Net alongside the corresponding ground truth. The pixel-wise accuracy metrics provide numerical insights, while the visual comparison showcases the model's capability in capturing semantic information.

4 Comparison: U-Net vs. CGAN

4.1 Accuracy Comparison

| Model | Pixel-wise Accuracy on Trainset (%) | Pixel-wise Accuracy on Validationset (%) |
|-------|-------------------------------------|--|
| U-Net | 56.75 | 59.88 |
| CGAN | 46.41 | 37.26 |

Table 3: Comparison of Pixel-wise Accuracy between U-Net and CGAN

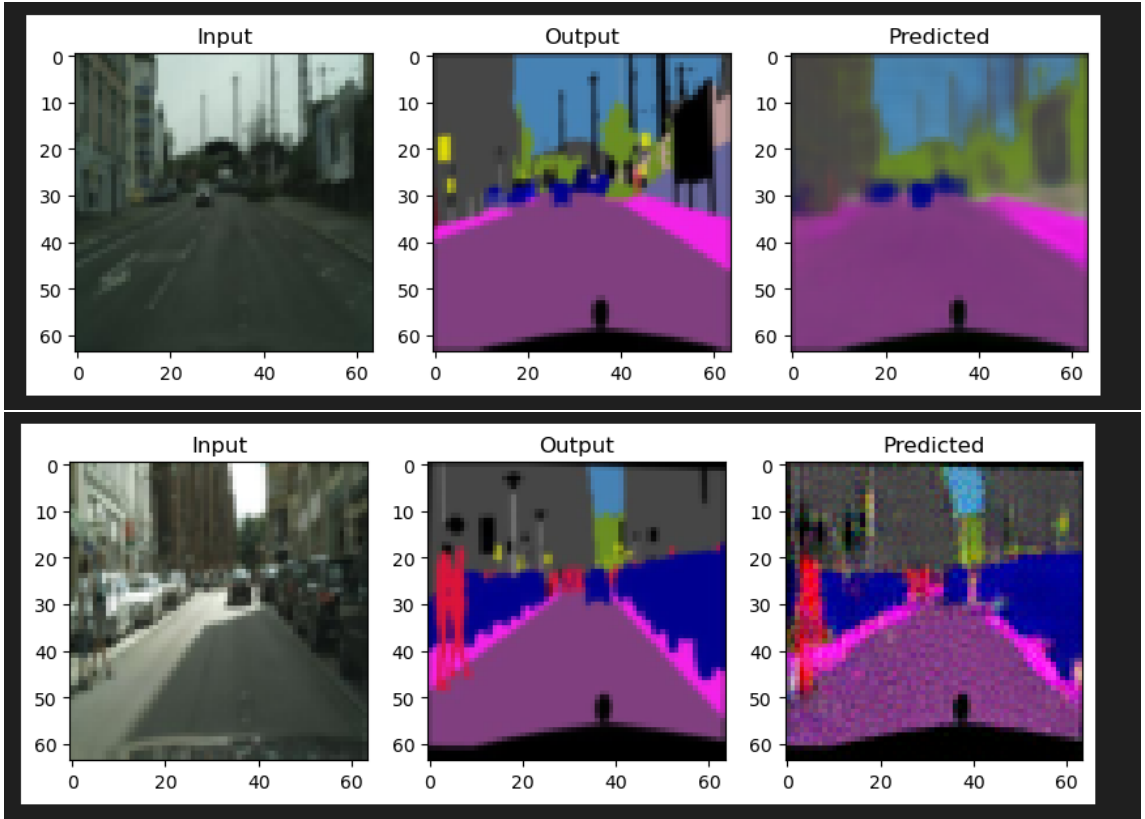


Figure 4: Comparison of Generated Images - U-Net (top) vs CGAN (bottom)

4.2 Generated Images Comparison

4.3 Comment

Even though the pixel-wise accuracy shows that the U-Net gives a better result, the observation reveals that the CGAN produces better-segmented regions with much clearer edges. However, CGAN segmented regions do contain lots of noises, which reduce the accuracy from the pixel-wise accuracy calculation method. Furthermore, the model result still shows to be converging after epochs but slowly. Due to a lack of time and resources, the training of CGAN stopped at 100 epochs. Hence, the CGAN might need more time for training to reach its peak performance.