

DLP  
Lab2: Binary Semantic Segmentation  
Report

林睿騰(Rui-Teng Lin)  
313551105

March 24, 2025

# Contents

1	Implementation Details	3
1.1	Training . . . . .	3
1.1.1	Initialize . . . . .	3
1.1.2	Train . . . . .	4
1.1.3	Evaluate . . . . .	4
1.2	Evaluation . . . . .	5
1.3	Inference . . . . .	5
1.4	Model Architecture . . . . .	6
1.4.1	UNet . . . . .	6
1.4.2	ResNet34Unet . . . . .	9
1.5	Loss Function . . . . .	12
2	Data Preprocessing	13
2.1	Preprocessing . . . . .	13
2.2	What makes the preprocessing method unique . . . . .	14
3	Analyze the experiment results	14
3.1	Training Curves . . . . .	14
3.2	Performance Metrics . . . . .	14
3.3	Visualization of Segmentation Results . . . . .	15
3.4	Discussion . . . . .	16
4	Execution steps	16
4.1	How to run the code . . . . .	16
5	Discussion	17
5.1	Alternative architectures . . . . .	17
5.2	Potential Research Topics . . . . .	18
6	References	18

# 1 Implementation Details

## 1.1 Training

The training procedure is just a ordinary training procedure in PyTorch. For the implementation of the training process, I follow the steps below:

1. Load the model and dataset
2. Define the optimizer, learning rate scheduler, and loss function
3. Train the model
4. Evaluate the model after each epoch

For the detailed implementation of each step:

### 1.1.1 Initialize

Load the model and train, valid dataset, and define the optimizer, learning rate scheduler, and loss function.

```
1 # Load the model
2 if args.model == "unet":
3     model = UNet()
4 elif args.model == "resnet34_unet":
5     model = ResNet34Unet()
6 else:
7     raise ValueError(f"Model {args.model} not found")
8
9 model.to(device)
10
11 # Load the dataset
12 train_dataset = load_dataset(args.data_path, "train")
13 valid_dataset = load_dataset(args.data_path, "valid")
14
15 train_loader = DataLoader(
16     train_dataset, batch_size=args.batch_size, shuffle=True,
17     num_workers=32
18 )
19 valid_loader = DataLoader(
20     valid_dataset, batch_size=args.batch_size, shuffle=False,
21     num_workers=32
22 )
23 optimizer = torch.optim.Adam(model.parameters(), lr=args.learning_rate)
24 scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer,
25     T_max=args.epochs)
26 loss_fn = torch.nn.BCELoss()
```

Listing 1: train.py: Initialize

### 1.1.2 Train

Train the model for a specified number of epochs, and use Binary Cross Entropy loss function and Dice loss function as the loss function.

```
1 for epoch in range(args.epochs):
2     model.train()
3     train_dice_score, train_bce_loss, train_dice_loss = 0, 0, 0
4
5     for batch in tqdm(train_loader, desc=f"Epoch {epoch +
6         1}/{args.epochs}"):
7         images = batch["image"].to(device).float()
8         masks = batch["mask"].to(device)
9
10        optimizer.zero_grad()
11        pred_masks = model(images)
12
13        d_score = dice_score(pred_masks, masks)
14        b_loss = loss_fn(pred_masks, masks)
15        d_loss = dice_loss(pred_masks, masks)
16
17        loss = b_loss + d_loss
18        loss.backward()
19        optimizer.step()
20
21        train_dice_score += d_score.item()
22        train_bce_loss += b_loss.item()
23        train_dice_loss += d_loss.item()
```

Listing 2: train.py: Train

### 1.1.3 Evaluate

Evaluate on validation set and save the best model (and log the results on wandb)

```
1 eval_dice_score, eval_bce_loss, eval_dice_loss = evaluate(
2     model, valid_loader, device
3 )
4 if eval_dice_score > best_eval_dice_score:
5     best_eval_dice_score = eval_dice_score
6     torch.save(
7         model.state_dict(),
8         f"saved_models/{args.model}_{epoch}_{eval_dice_score:.4f}.pth",
9     )
10
11 if args.wandb:
12     wandb.log(
13         {
14             "train/dice_score": train_dice_score / len(train_loader),
15             "train/bce_loss": train_bce_loss / len(train_loader),
16             "train/dice_loss": train_dice_loss / len(train_loader),
17             "valid/dice_score": eval_dice_score,
18             "valid/bce_loss": eval_bce_loss,
19             "valid/dice_loss": eval_dice_loss,
20         }
```

Listing 3: train.py: Evaluate

## 1.2 Evaluation

For evaluation, I use the same procedure as the training process, the only difference is the dataset and the model is not updated.

```

1 def evaluate(net, data, device):
2     # implement the evaluation function here
3     net.eval()
4     dice_scores = []
5     bce_losses = []
6     dice_losses = []
7     bce_loss = torch.nn.BCELoss()
8     with torch.no_grad():
9         for batch in data:
10             images = batch["image"].to(device).float()
11             masks = batch["mask"].to(device)
12             pred_masks = net(images)
13             dice_scores.append(dice_score(pred_masks, masks).item())
14             bce_losses.append(bce_loss(pred_masks, masks).item())
15             dice_losses.append(dice_loss(pred_masks, masks).item())
16     return np.mean(dice_scores), np.mean(bce_losses), np.mean(dice_losses)

```

Listing 4: evaluate.py: Evaluate

## 1.3 Inference

For inference, I use the same procedure as the evaluation process, and this time, I need to load the state dict of the model and calculate the accuracy on the test set.

```

1 def inference(model, device):
2     if "resnet34" in args.model:
3         model = ResNet34Unet(in_channels=3)
4     else:
5         model = UNet(in_channels=3)
6
7     model.load_state_dict(torch.load(args.model))
8     model.to(device)
9     model.eval()
10
11     test_dataset = load_dataset(args.data_path, "test")
12     test_loader = DataLoader(test_dataset, batch_size=args.batch_size,
13                             shuffle=False)
14
15     dice_score, _, _ = evaluate(model, test_loader, device=device)
16     print(f"Dice score: {dice_score:.4f}")

```

Listing 5: inference.py: Inference

## 1.4 Model Architecture

### 1.4.1 UNet

For the UNet model, I use almost the same architecture as the one in the paper. The only difference is that in the original paper, the DownConv layers do not use padding, so the output size is smaller than the input size divided by 2. However, in this implementation, I use padding, so the output size is the same as the input size divided by 2, since by this way, it is easier to implement the upsampling layer, and also easier to integrate with the ResNet34.

The architecture of the UNet is shown below (the size of the input is 256x256, the size and the dimension of each layer output is annotated in the diagram):

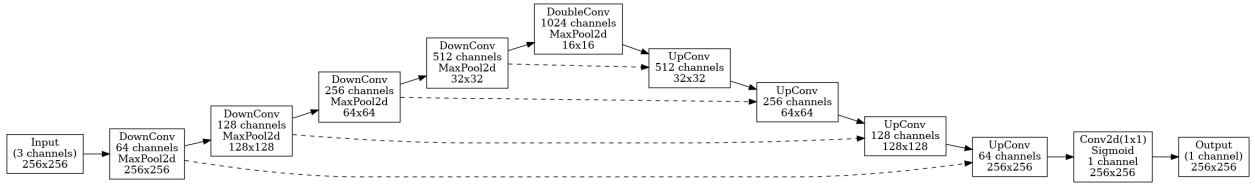


Figure 1: UNet Architecture

For the implementation of the UNet, I refer to the code from a public repository milesial/Pytorch-UNet.

To simply describe the architecture, it can be divided into 3 parts:

1. DoubleConv layers: The double convolutional layers, which is a series of convolutional layers with batch normalization and ReLU activation functions.
2. DownConv layers: The downsampling path of the UNet, which is a DoubleConv layer with max-pooling layers.
3. UpConv layers: The upsampling path of the UNet, which is a DoubleConv layer with bilinear upsampling layers.
4. Middle layer: The middle layer of the UNet, which is a DoubleConv layer.

**DoubleConv** The DoubleConv layer is a series of convolutional layers with batch normalization and ReLU activation functions. There was not BatchNorm in the original paper, but I add it for better performance.

```
1 class DoubleConv(nn.Module):
2     def __init__(self, in_channels, out_channels):
3         super(DoubleConv, self).__init__()
4
5         layers = [
6             nn.Conv2d(in_channels, out_channels, kernel_size=3, padding=1,
7                       bias=False),
8             nn.BatchNorm2d(out_channels),
9             nn.ReLU(inplace=True),
10            nn.Conv2d(out_channels, out_channels, kernel_size=3, padding=1,
11                     bias=False),
12            nn.BatchNorm2d(out_channels),
13            nn.ReLU(inplace=True),
```

```

12     ]
13
14     self.nn = nn.Sequential(*layers)
15
16 def forward(self, x):
17     return self.nn(x)

```

Listing 6: models/unet.py: DoubleConv

**DownConv** The DownConv layer is a DoubleConv layer with max-pooling layers. It is used to reduce the spatial size of the feature map and expand the channel dimension.

```

1 class DownConv(nn.Module):
2 def __init__(self, in_channels, out_channels):
3     super(DownConv, self).__init__()
4
5     layers = [nn.MaxPool2d(kernel_size=2), DoubleConv(in_channels,
6                                                         out_channels)]
7
8     self.nn = nn.Sequential(*layers)
9
10 def forward(self, x):
11     return self.nn(x)

```

Listing 7: models/unet.py: DownConv

**UpConv** The UpConv layer is a DoubleConv layer with bilinear upsampling layers. It is used to increase the spatial size of the feature map and reduce the channel dimension. The input from previous layer (the mid / UpConv layer) is concatenated with the output from the DownConv layer at the same spatial size after upsampling, and then the DoubleConv layer is applied.

```

1 class UpConv(nn.Module):
2 def __init__(self, in_channels, out_channels):
3     super(UpConv, self).__init__()
4
5     self.up = nn.ConvTranspose2d(in_channels, out_channels,
6                                   kernel_size=2, stride=2)
7     self.conv = DoubleConv(in_channels, out_channels)
8
9 def forward(self, x1, x2):
10     x = self.up(x1)
11     x = torch.cat([x, x2], dim=1)
12     return self.conv(x)

```

Listing 8: models/unet.py: UpConv

**Middle** The middle layer is a DoubleConv layer, which is used to connect the upsampling path and the downsampling path.

```

1 self.mid = DoubleConv(down_channels[-1], up_channels[0])

```

Listing 9: models/unet.py: UNet

UNet The UNet model is a combination of the DoubleConv, DownConv, UpConv, and Middle layers. The output activation function is sigmoid, so we can get the binary segmentation mask and use Binary Cross Entropy loss function. I think the most confusing part is the skip connection, which is used to connect the output from the DownConv layer and the input from the UpConv layer at the same spatial size. Make sure which one should be popped from the stack is important.

```

1  class UNet(nn.Module):
2  def __init__(
3      self,
4      in_channels: int = 3,
5      down_channels: list[int] = [64, 128, 256, 512],
6      up_channels: list[int] = [1024, 512, 256, 128, 64],
7      out_channels: int = 1,
8  ):
9      super(UNet, self).__init__()
10
11     self.in_conv = DoubleConv(in_channels, down_channels[0])
12
13     self.down = nn.ModuleList()
14     self.up = nn.ModuleList()
15
16     for i in range(len(down_channels) - 1):
17         self.down.append(DownConv(down_channels[i], down_channels[i +
18                                     1]))
19
20     self.mid = DownConv(down_channels[-1], up_channels[0])
21
22     for i in range(len(up_channels) - 1):
23         self.up.append(UpConv(up_channels[i], up_channels[i + 1]))
24
25     self.out = nn.Sequential(
26         nn.Conv2d(up_channels[-1], out_channels, kernel_size=1),
27         nn.Sigmoid(),
28     )
29
30 def forward(self, x):
31     x = self.in_conv(x)
32
33     x_rec = [x]
34     for down in self.down:
35         x_rec.append(down(x))
36         x = x_rec[-1]
37
38     x = self.mid(x)
39
40     for up in self.up:
41         x = up(x, x_rec.pop())
42     return self.out(x)

```

Listing 10: models/unet.py: UNet



### 1.4.2 ResNet34Unet

For the ResNet34Unet model, I use the ResNet34 as the encoder and the UNet as the decoder. The ResNet34 model architecture is mostly same as the original paper and also the image in Lab2 Spec.

The architecture of the ResNet34Unet is shown below (the size of the input is 256x256, the size and the dimension of each layer output is annotated in the diagram):

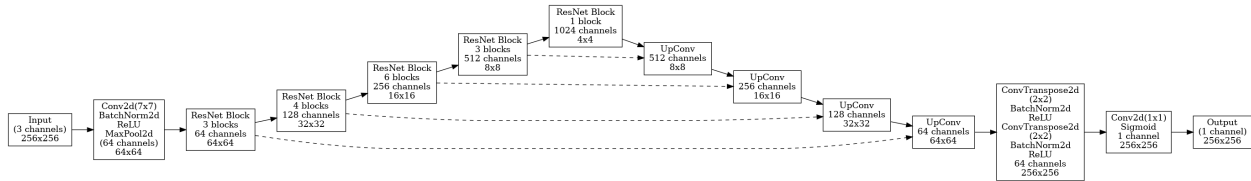


Figure 2: ResNet34Unet Architecture

For the left part of the ResNet34Unet, it is the same as the original ResNet34 model. And for the right part, it is the UNet model with two extra ConvTranspose2d layers to upsample the feature map to the original size (since the input is lowered by 4 times at the first layer of the ResNet34, so we need to upsample it by 4 times to get the original size).

To break down the ResNet34Unet, it can be divided into 3 parts:

1. ResBlock: The Residual Block of the ResNet34, used to extract the features from the input.
2. UpConv: The upsampling path of the UNet, which is a ConvTranspose2d layer with bilinear upsampling layers, used to upsample the feature map to the original size.

**ResBlock** The ResBlock layer is the same as the Residual Block of the ResNet34 model.

There is a classmethod `make_layer` to make a series of ResBlock layers, which is used to make a series of ResBlock layers, and automatically add the downsampling path to the first ResBlock layer when the channel number is doubled.

```

1 class ResBlock(nn.Module):
2     def __init__(self, in_channels: int, out_channels: int, down_sample:
3         bool = False):
4         super(ResBlock, self).__init__()
5
6         self.residual = nn.Sequential(
7             nn.Conv2d(
8                 in_channels,
9                 out_channels,
10                kernel_size=3,
11                padding=1,
12                stride=2 if down_sample else 1,
13                bias=False,
14            ),
15            nn.BatchNorm2d(out_channels),
16            nn.ReLU(inplace=True),
17            nn.Conv2d(out_channels, out_channels, kernel_size=3,
18                padding=1, bias=False),
19            nn.BatchNorm2d(out_channels),
20        )

```

```

19         self.shortcut = (
20             nn.Sequential(
21                 nn.Conv2d(
22                     in_channels,
23                     out_channels,
24                     kernel_size=1,
25                     bias=False,
26                     stride=2 if down_sample else 1,
27                 ),
28                 nn.BatchNorm2d(out_channels),
29             )
30         )
31         if down_sample
32         else nn.Identity()
33     )
34
35     self.relu = nn.ReLU(inplace=True)
36
37     def forward(self, x):
38         return self.relu(self.residual(x) + self.shortcut(x))
39
40     @classmethod
41     def make_layer(
42         cls, in_channels: int, out_channels: int, blocks: int,
43         down_sample: bool = False
44     ):
45         layers = [cls(in_channels, out_channels, down_sample)]
46         for _ in range(1, blocks):
47             layers.append(cls(out_channels, out_channels))
48         return nn.Sequential(*layers)

```

Listing 11: models/resnet34\_unet.py: ResBlock

UpConv Same as one in the UNet model (see 8).

**ResNet34Unet** The ResNet34Unet model is a combination of the ResBlock and UpConv layers. The output activation function is sigmoid, so we can get the binary segmentation mask and use Binary Cross Entropy loss function.

The implementation is different from the one in the Lab2 Spec, since the original UNet concatenates tensors with same dimensions, so I think we should follow this rule. To meet this requirement, I use another ResBlock as the bottleneck, then forward the output with the previous ResBlocks' output to the UNet decoder.

By following this rule, the input of the first UpConv layer is by default 1024 dim (bottleneck output)  $-(\text{ConvTranspose2d}) \rightarrow 512 \text{ dim} + 512 \text{ dim}$  (last ResBlock output) = 1024 dim, and undergoes a DoubleConv layer to reduce the channel dimension to 512 dim.

So, according to the above description, the implementation is as follows:

```

1 class ResNet34Unet(nn.Module):
2     def __init__(
3         self,
4         in_channels: int = 3,
5         out_channels: int = 1,

```

```

6     blocks: list[int] = [3, 4, 6, 3],
7     down_channels: list[int] = [64, 64, 128, 256, 512],
8     up_channels: list[int] = [1024, 512, 256, 128, 64],
9 ):
10     super(ResNet34Unet, self).__init__()
11
12     self.in_conv = nn.Sequential(
13         nn.Conv2d(
14             in_channels,
15             down_channels[0],
16             kernel_size=7,
17             padding=3,
18             stride=2,
19             bias=False,
20         ),
21         nn.BatchNorm2d(down_channels[0]),
22         nn.ReLU(inplace=True),
23         nn.MaxPool2d(kernel_size=3, stride=2, padding=1),
24     )
25
26     self.down = nn.ModuleList()
27     for i in range(len(down_channels) - 1):
28         self.down.append(
29             ResBlock.make_layer(
30                 down_channels[i],
31                 down_channels[i + 1],
32                 blocks[i],
33                 down_sample=(down_channels[i] != down_channels[i +
34                     1])),
35         )
36
37     self.mid = ResBlock.make_layer(
38         down_channels[-1],
39         up_channels[0],
40         1,
41         down_sample=(down_channels[-1] != up_channels[0]),
42     )
43
44     self.up = nn.ModuleList()
45     for i in range(len(up_channels) - 1):
46         self.up.append(UpConv(up_channels[i], up_channels[i + 1]))
47
48     self.out = nn.Sequential(
49         nn.ConvTranspose2d(
50             up_channels[-1], up_channels[-1], kernel_size=2, stride=2
51         ),
52         nn.BatchNorm2d(up_channels[-1]),
53         nn.ReLU(inplace=True),
54         nn.ConvTranspose2d(
55             up_channels[-1], up_channels[-1], kernel_size=2, stride=2
56         ),
57         nn.BatchNorm2d(up_channels[-1]),
58         nn.ReLU(inplace=True),

```

```

59         nn.Conv2d(up_channels[-1], out_channels, kernel_size=1),
60         nn.Sigmoid(),
61     )
62
63     def forward(self, x):
64         x = self.in_conv(x)
65         x_rec = []
66         for layer in self.down:
67             x = layer(x)
68             x_rec.append(x)
69         x = self.mid(x)
70         for layer in self.up:
71             x = layer(x, x_rec.pop())
72         return self.out(x)

```

Listing 12: models/resnet34\_unet.py: ResNet34Unet

## 1.5 Loss Function

For loss function, two loss functions are used: Binary Cross Entropy loss function and Dice loss function.

**Binary Cross Entropy loss function** The Binary Cross Entropy loss function is a loss function that measures the performance of a binary classification model, so it is suitable for the binary segmentation task. Since it is implemented in PyTorch, so we won't need to dig into the details of the implementation.

**Dice loss function** The Dice loss function is a loss function that measures the performance of a binary segmentation model, so it is suitable for the binary segmentation task.

$$\text{DiceLoss}(p, y) = 1 - \frac{2 \sum_{i=1}^N y_i p_i}{\sum_{i=1}^N y_i + \sum_{i=1}^N p_i} \quad (1)$$

and surprisingly, the gradient of the DiceLoss can be represented as a simple formula (in the referenced repository):

$$\frac{\partial \text{DiceLoss}(p, y)}{\partial p} = 1 - \text{DiceLoss}(p, y) \quad (2)$$

after we figure out the formula of the loss and its gradient, we can implement it in the code.

```

1 def dice_score(pred_mask, gt_mask):
2     # implement the Dice score here
3     assert pred_mask.shape == gt_mask.shape
4
5     if pred_mask.ndim == 3: # expand (C, H, W) -> (1, C, H, W)
6         pred_mask = pred_mask.unsqueeze(1)
7
8     pred_mask = torch.where(pred_mask > 0.5, True, False)
9     gt_mask = torch.where(gt_mask > 0.5, True, False)
10
11     common = torch.sum(pred_mask & gt_mask, dim=(1, 2, 3)) # common
    pixels between pred and gt

```

```

12     union = torch.sum(pred_mask, dim=(1, 2, 3)) + torch.sum(gt_mask,
13         dim=(1, 2, 3))
14     union = torch.where(union == 0, 1, union) # avoid division by zero
15
16     return (2 * common / union).mean() # average over the batch
17
18
19 def dice_loss(pred_mask, gt_mask):
20     return 1 - dice_score(pred_mask, gt_mask)

```

Listing 13: utils.py: dice\_loss

## 2 Data Preprocessing

### 2.1 Preprocessing

For preprocessing, I use the following steps:

1. Resize the image to 256x256
2. Randomly crop the image to 256x256 or rotate the image by at most 30 degrees
3. Randomly flip the image horizontally or vertically
4. Randomly adjust the brightness or gamma of the image
5. Randomly add Gaussian noise or blur to the image
6. Normalize the image with the mean and standard deviation of the ImageNet dataset

and implement it in the `load_dataset` function using the functions from `albumentations` library.

```

1 A.Compose([
2     A.Resize(256, 256),
3     A.OneOf([
4         A.RandomResizedCrop((256, 256), scale=(0.8, 1.0), ratio=(0.75,
5             1.33), p=1.0),
6         A.Rotate(limit=30, p=1.0),
7     ], p=0.5),
8     A.OneOf([
9         A.HorizontalFlip(p=1.0),
10        A.VerticalFlip(p=1.0),
11    ], p=0.5),
12    A.OneOf([
13        A.RandomBrightnessContrast(p=1.0),
14        A.RandomGamma(p=1.0),
15    ], p=0.3),
16    A.OneOf([
17        A.GaussNoise(p=1.0),
18        A.GaussianBlur(p=1.0),
19    ], p=0.2),
20    A.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
    ToTensorV2(),

```

## 2.2 What makes the preprocessing method unique

I choose some transformations of geometric, photometric, and noise to augment the data. But if we impose too many transformations, the image would be too different from the original image, and the model would not be able to learn the features.

So, I use the `OneOf` function to choose one of the transformations to apply to the image. And use the parameter  $p$  to control the probability of each transformation, prevent the image from being too different.

By this method, the model can learn from moderately augmented data, and the performance is better than using only the original data.

## 3 Analyze the experiment results

### 3.1 Training Curves

The training curves are shown in the following figure 3.

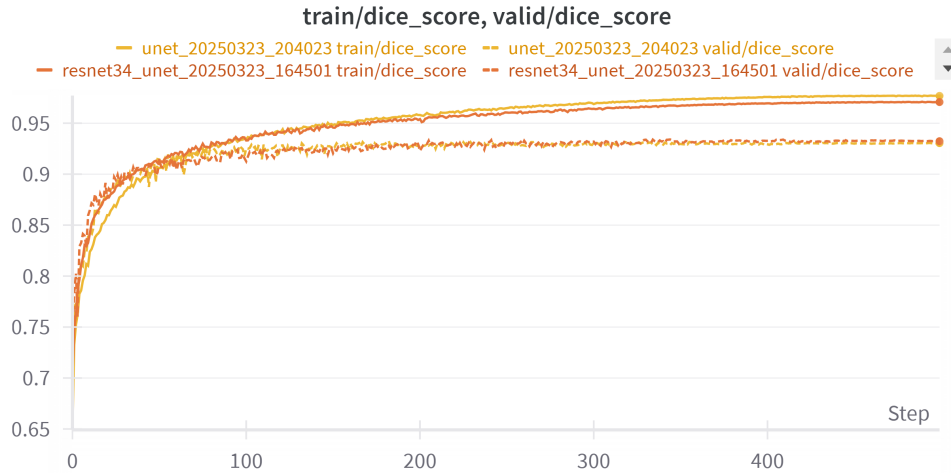


Figure 3: Training curves: Dice score. X-axis is the number of epochs, Y-axis is the Dice score.

There two pairs of curves, one is for ResNet34 UNet (Orange) and the other is for UNet (Yellow). And the solid line is the dice score of the training set, and the dashed line is the dice score of the validation set. We can see from the figure, the performance of the two models are almost the same, but the ResNet34 UNet converges faster than UNet, and UNet finally converges to a higher train accuracy than ResNet34 UNet. But the validation accuracy of ResNet34 UNet is higher than UNet, which means that ResNet34 UNet is less likely to overfit than UNet.

### 3.2 Performance Metrics

I use the Dice score on test set to evaluate the performance of the model. With the parameter given in the execution section 17, the performance of the model is shown in the following table 1

and figure 4.

Model	Dice score
UNet	0.9318
ResNet34 UNet	0.9330

Table 1: Dice score of the test set. The above is the ResNet34 UNet, and the below is the UNet.

```

(Lab2-3.12) rtlin@admin7428-SYSTEM-PRODUCT-NAME:~/2025_Spring_Deep-Learning-Labs/Lab2$ python src/inference.py --model saved_models/unet_0.9318.pth --data_path dataset/
Dice score: 0.9318
(Lab2-3.12) rtlin@admin7428-SYSTEM-PRODUCT-NAME:~/2025_Spring_Deep-Learning-Labs/Lab2$ python src/inference.py --model saved_models/resnet34_unet_0.9330.pth --data_path d
ataset/
Dice score: 0.9330

```

Figure 4: Dice score of the test set. The above is the UNet, and the below is the ResNet34 UNet.

### 3.3 Visualization of Segmentation Results

The visualization of the segmentation results is shown in the following figure. There are five samples in the figure, the first row is the original image, the second row is the ground truth mask, the third row is the prediction from the UNet, and the fourth row is the prediction from the ResNet34 UNet.

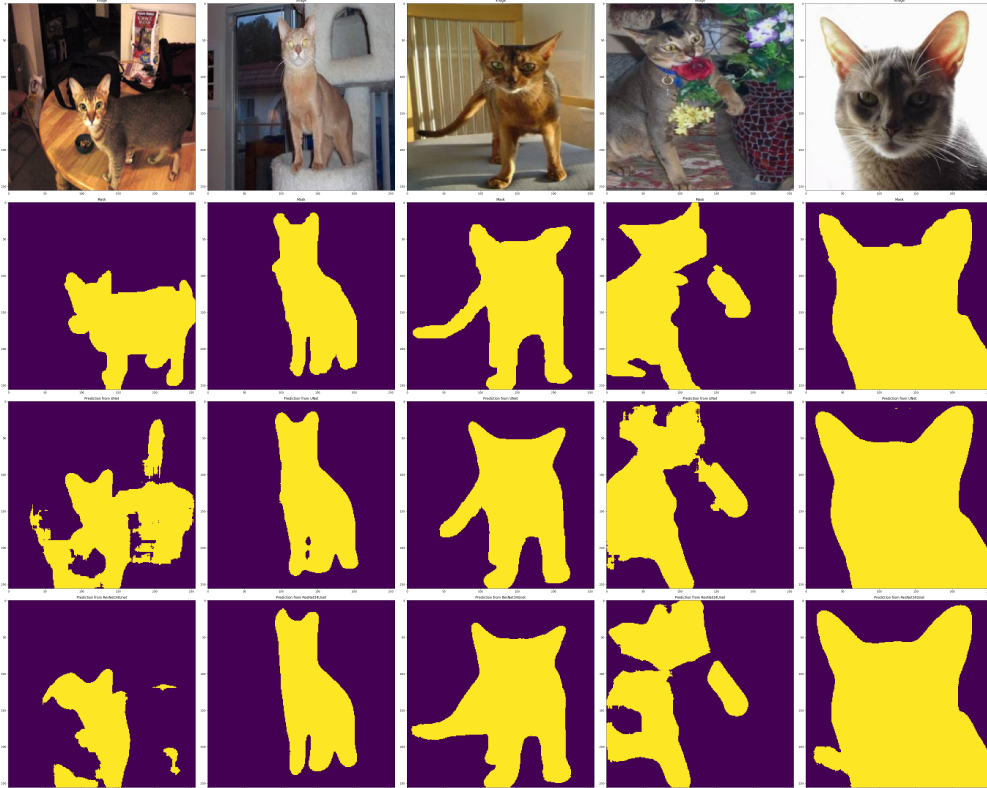


Figure 5: Visualization of segmentation results. The above is the UNet, and the below is the ResNet34 UNet.

We can see from the figure, Resnet34 UNet has a better performance than UNet, it can predict the small objects better, and it also performs when dealing with the boundary of the objects. This

behavior is more obvious in sample 1 and 3, where the prediction from Unet is totally distorted, but the prediction from Resnet34 UNet is still good.

### 3.4 Discussion

Throughout this project, I firstly experiment with both of the two models, but I found that ResNet34 UNet has a 3x faster training speed than UNet, so I use ResNet34 UNet for following experiments.

Secondly, I conducted experiments on learning rate, batch size, and number of epochs, and I found that the learning rate have a significant impact on the performance of the model, but the batch size and number of epochs have a little impact. For example, when the learning rate is set to  $1e-4$  to  $1e-3$ , the performance of ResNet34 UNet about the same, as for learning rate  $1e-5$ , it converges too slow that it cant converge to a good accuracy in 500 epochs. The batch size is set to 64, is not because it leads to a better accuracy, but because it leads to A LOT faster training speed without out of memory error. And the number of epochs is set to 500, is because I found that the performance of the model is almost saturated after 500 epochs, and it is not necessary to train the model for more epochs.

After I decided the parameters, I conducted experiments on the transformation of the images, and I found that the performance of the model is not sensitive to the transformation of the images. Originally, I conducted lots of transformations on the training image, and the performance is almost the same as those experiment with half of the transformations. After then, I found that normalization and rotation have a significant impact on the performance of the model, and other transformations have little impact on the performance of the model. So, as a result, I kept some of the transformations, and using OneOf to prevent over-transformation, and normalize the image at the end.

## 4 Execution steps

### 4.1 How to run the code

Before running the code, make sure you are in the right directory.

```
$ tree
```

```
|—— dataset
|—— saved_models
|—— src
```

You can run the code by executing the following command:

First, install the dependencies:

```
1 pip install -r requirements.txt
```

Listing 15: Install the dependencies

Then, run the training code:

```
1 # Run the code for training
2 python src/train.py --data_path dataset/oxford-iiit-pet/
```

Listing 16: How to run the training code

there are other parameters you can use to customize the training process:



- `--device`: the device to run the code on (default: `cuda:1`)
- `--seed`: the seed for the random number generator (default: 88)
- `--epochs`: the number of epochs to train the model (default: 500)
- `--batch_size`: the batch size for the training process (default: 64)
- `--learning_rate`: the learning rate for the training process (default: `7e-4`)
- `--model`: the model to use for the training process (default: `resnet34_unet`), you can choose from: `resnet34_unet`, `unet`
- `--wandb`: whether to use wandb for logging (default: `True`)

For my own results, I used the following commands:

```

1 # ResNet34 UNet
2 python src/train.py \
3     --data_path dataset/oxford-iiit-pet/ \
4     -d cuda:1 \
5     -s 228922
6
7 # UNet
8 python src/train.py \
9     --data_path dataset/oxford-iiit-pet/ \
10    -d cuda:1 \
11    -s 228922 \
12    -m unet

```

Listing 17: How to run the training code for my own results

After training, you can run the inference code:

```

1 # Run the code for inference
2 # Replace the model_path with the path to the model you want to use for
  the inference
3 # Model name should contain resnet34_unet or unet
4 python src/inference.py --data_path dataset/oxford-iiit-pet/ --model_path
  saved_models/model.pth

```

Listing 18: How to run the inference code

## 5 Discussion

### 5.1 Alternative architectures

I think in this task, we can use more complex architectures to improve the performance of the model, since the input image contains various conditions, including different brightness, different contrast, different size, different angle, and the object itself can have a weird shape (especially cats).

So, as an alternative, we can conduct experiments on the most famous model: Transformer. We can leverage its powerful attention mechanism to improve the performance of the model, let it more focus on the important parts of the image. And more on, we can use Mixture-of-Experts to improve the performance of the model, let it more robust to different conditions.

## 5.2 Potential Research Topics

Using transformer is a good idea, but everyone is using it nowadays, and the parameters is getting larger and larger. So, I think it is a good idea to distill the trained transformer model on this task into a smaller model, trying to keep the performance of the model and reduce the parameters. Or, on the other hand, we can try to use a small model, but try some fancy method to deal with the input image with various conditions, finally get comparable performance with transformer-based models.

## 6 References

milesial/Pytorch-UNet