

# NYCU DL Lab3 Binary Semantic Segmentation

313551073 顏琦恩

## I. Overview of your lab 3

In this lab, I implemented UNet and ResNet34-UNet from the paper U-Net: Convolutional Networks for Biomedical Image Segmentation and Deep learning-based pelvic levator hiatus segmentation from ultrasound images.

In the training section, I'd trained two models, UNet and ResNet34-UNet. I use the dataset provided by TAs, which is The Oxford-IIIT Pet Dataset, to train my model.

Moreover, I'd implemented some tool functions in file utils.py, which are dice score calculating function, three plotting functions, functions for loading and saving models and loss. Plotting functions are for the loss curve and the predicted masks. The saving function for loss is used when one needs to keep on training the same model. To draw a complete loss curve, this function saves losses into a txt file, one can use the loading function to read the txt file and draw the complete loss curve.

```
# 畫多張圖、ground truth、predicted mask，一行最多8張圖
def plot_grid(image, mask, pred_mask):
    fig = plt.figure(figsize=(10, 12))
    fig.suptitle('Comparison of Original Image, Mask, and Predicted Mask', fontsize=12)

    fig.add_subplot(3, 1, 1)
    plt.imshow(T.ToPILImage()(torchvision.utils.make_grid(image / 255, nrow=8)))
    plt.axis('off')
    plt.title("Original Image")

    fig.add_subplot(3, 1, 2)
    plt.imshow(T.ToPILImage()(torchvision.utils.make_grid(mask, nrow=8)))
    plt.axis('off')
    plt.title("Ground Truth Masks")

    fig.add_subplot(3, 1, 3)
    plt.imshow(T.ToPILImage()(torchvision.utils.make_grid(pred_mask, nrow=8)))
    plt.axis('off')
    plt.title("Predicted Masks")

    plt.show()
```

**Fig. 1 Function for plotting multiple images at the same time**

## II. Implementation Details

### A. Details of your training, evaluating, inferencing code

In my training section, I applied binary cross-entropy as my loss function and Adam as my optimizer. I used binary cross-entropy because this lab only contains two classes: foreground and background. To avoid the error of cuda out of memory, I clean up the gpu

cache after finishing every epoch. Also, to keep training the model, I added a variable `retrain_model` to keep on training my previous trained model and save losses of training to draw the complete training loss curve. During training, I check if the dice score of the validation dataset gets better, and save the model which has the highest dice score in the process. At the end of the training function, I plot the loss curve of the training section.

```

14 def train(args):
15     # implement the training function here
16     train_dataloader = DataLoader(opData.load_dataset(args.data_path, 'train'), batch_size=args.batch_size, shuffle=True)
17     valid_dataloader = DataLoader(opData.load_dataset(args.data_path, 'valid'), batch_size=args.batch_size, shuffle=True)
18
19     if args.model == 'unet':
20         model = unet.UNet(3).to(args.device)
21     elif args.model == 'resnet':
22         model = resnet34_unet.ResNet34_UNet(3).to(args.device)
23
24     criterion = nn.BCELoss() # 用於只有binary class的cross entropy
25     optimizer = optim.Adam(model.parameters(), lr=args.learning_rate)
26
27     losses = []
28     scores = []
29     max_score = 0
30
31     # 如果要繼續train同一個model, 可以使用這個部分, loss會寫入txt檔中, 所以可以畫完整的loss curve
32     if args.retrain_model != '':
33         model, optimizer, max_score = load_model(f'./saved_models/{args.retrain_model}', model, optimizer)
34         losses = load_loss(args.model)
35
36     # start training
37     last_epoch = len(losses) # 從上次訓練到的epoch繼續畫loss curve, 如果是train from scratch, 就從epoch 0開始畫loss curve
38
39     model.train()
40     for epoch in range(args.epochs):
41
42         epoch_cost = []
43         epoch_score = []
44         for sample in tqdm(train_dataloader):
45             img = torch.Tensor(sample['image']).float().to(args.device)
46             mask = torch.Tensor(sample['mask']).float().to(args.device)
47             optimizer.zero_grad()
48
49             pred = model(img)
50             loss = criterion(pred, mask)
51             loss.backward()
52             optimizer.step()
53
54             epoch_cost.append(loss)
55             epoch_score.append(dice_score(pred, mask))
56
57             losses.append((sum(epoch_cost)/len(epoch_cost)).item())
58             scores.append((sum(epoch_score)/len(epoch_score)).item())
59             print(f'[Epoch {epoch+1}] loss = {losses[-1]:.9f}, average dice score = {scores[-1]:.9f}')
60
61         # validation
62         valid_score = evaluate(model, valid_dataloader, args.device)
63         print(f'Validation dice score = {valid_score}')
64         if valid_score > max_score:
65             save_model(f'{args.save_path + args.model}.pth', model, optimizer, epoch, valid_score)
66             max_score = valid_score
67         save_loss(args.model, losses[last_epoch+epoch:]) # 將loss寫入txt檔
68         # 避免gpu out of memory
69         gc.collect()
70         torch.cuda.empty_cache()
71
72     plot_loss(losses, args.model)

```

**Fig. 2 Details of training code**

The structure of evaluating and inferencing code is familiar with each other. The difference between them is that I need to load my trained model for inference, but don't need to load the model in the evaluating section, since I always pass the model to evaluate into the evaluating function.

```

5  def evaluate(net, data, device):
6      # implement the evaluation function here
7      avg_score = []
8      net.eval()
9      with torch.no_grad():
10         for sample in tqdm(data):
11             img = torch.tensor(sample['image']).float().to(device)
12             mask = torch.tensor(sample['mask']).float().to(device)
13
14             pred = net(img)
15             avg_score.append(dice_score(pred, mask))
16
17     return sum(avg_score)/len(avg_score)

```

**Fig. 3 Details of evaluating code**

```

20  if __name__ == '__main__':
21      args = get_args()
22
23      dataloader = DataLoader(opData.load_dataset(args.data_path, 'test'), batch_size=args.batch_size, shuffle=True)
24      if args.model_to_use == 'unet':
25          model = load_model(f'./saved_models/{args.model}', unet.UNet(3))
26      elif args.model_to_use == 'resnet':
27          model = load_model(f'./saved_models/{args.model}', resnet34_unet.ResNet34_UNet(3))
28
29      avg_score = []
30      model.to(args.device)
31      model.eval()
32      with torch.no_grad():
33         for sample in tqdm(dataloader):
34             img = torch.tensor(sample['image']).float().to(args.device)
35             mask = torch.tensor(sample['mask']).float().to(args.device)
36
37             pred = model(img)
38             score = dice_score(pred, mask)
39             avg_score.append(score)
40             if max(avg_score) == score and args.plot == 'y':
41                 plot_grid(img, mask, pred)
42     print(f'Inference dice score = {sum(avg_score)/len(avg_score)}')

```

**Fig. 4 Details of inferencing code**

## B. Details of your model (UNet & ResNet34\_UNet)

The UNet is composed of downsample blocks and upsample blocks. Each downsample block contains a maxpool, two 2d convolutions, relus, and batch normalizations, while the upsample block contains a transposed convolution and two 2d convolutions.

In the downsampling section, poolings are used to capture more abstract and higher resolution data, since they divide features into regions and keep the maximum value. For the kernel size, stride, and padding of convolutions and maxpool, I use the same value as the paper mentioned.

For the upsampling method, I used the transposed convolution, which is used to recover the spatial resolution of the image, the kernel size and stride is also used as the value mentioned in the paper.

```

# 透過downsample降低spatial resolution，來提取提取high-resolution，low-level的特徵
class DownsampleBlock(nn.Module):
    def __init__(self, in_channels, out_channels) -> None:
        super(DownsampleBlock, self).__init__()

        self.conv = nn.Sequential(
            nn.MaxPool2d(kernel_size=2, stride=2),
            Conv2dBlock(in_channels, out_channels, 3)
        )

    def forward(self, x):
        x = self.conv(x)
        return x

# 恢復spatial resolution，同時skip connection可以維持downsample時提取的特徵
class UpsampleBlock(nn.Module):
    def __init__(self, in_channels, out_channels) -> None:
        super(UpsampleBlock, self).__init__()

        self.upconv = nn.ConvTranspose2d(in_channels, out_channels, kernel_size=2, stride=2)
        self.conv = Conv2dBlock(in_channels, out_channels, kernel_size=3)

    def forward(self, x, prev_x):
        x = self.upconv(x)
        x = torch.cat([prev_x, x], dim=1) # skip connection

        x = self.conv(x)
        return x

```

**Fig. 5 The downsample and upsample block for UNet**

My ResNet34\_UNet took ResNet as the downsample part, which is the so-called encoder, while UNet did the upsample part, the decoder. The convolution blocks in the encoder are set to the same kernel size, stride, and padding as the ResNet paper.

In the upsample part, the channels of each decoder seem to be wrong, so I adjusted them to the same values as my Unet. Also, the paper of ResNet34-Unet mentioned that they add CBAM in each decoder, so I added it in my network, too. The remaining part of the decoder shared the same structure with the upsample part of UNet.

```

class DecoderBlock(nn.Module):
    def __init__(self, in_channels, out_channels, conv_in_channels=0) -> None:
        super(DecoderBlock, self).__init__()
        if conv_in_channels == 0:
            conv_in_channels = in_channels

        self.upconv = nn.ConvTranspose2d(in_channels, out_channels, kernel_size=2, stride=2, bias=False)
        self.unetconv = UnetConv(conv_in_channels, out_channels)
        self.cbam = CBAM(out_channels)

    def forward(self, x, prev_x):
        x = self.upconv(x)
        x = torch.cat([prev_x, x], dim=1)

        x = self.unetconv(x)
        x = self.cbam(x)
        return x

```



```

class EncoderBlock(nn.Module):
    def __init__(self, in_channels, out_channels, total_blocks, is_downsample=True) -> None:
        super(EncoderBlock, self).__init__()

        blocks = []
        blocks.append(ResBlock(in_channels, out_channels, kernel_size=3,
                               stride=1+1*is_downsample, is_downsample=is_downsample))
        for _ in range(1, total_blocks):
            blocks.append(ResBlock(out_channels, out_channels, kernel_size=3, stride=1))

        self.layer = nn.Sequential(*blocks)

    def forward(self, x):
        x = self.layer(x)
        return x

```

**Fig. 6 The decoder and encoder block of ResNet34\_UNet**

### III. Data Preprocessing

#### A. How you preprocessed your data?

I applied a random horizontal flip to the dataset with a probability of 0.5. The following figure is my flipping function.

```

def transform(image, mask, trimap):
    if np.random.rand() < 0.5:
        image = np.fliplr(image).copy()
        mask = np.fliplr(mask).copy()
        trimap = np.fliplr(trimap).copy()
    return {'image': image, 'mask': mask, 'trimap': trimap}

```

**Fig. 7 The preprocess function**

#### B. What makes your method unique?

Although just flipping images horizontally is a quite common way of data augmentation, it is an appropriate method in this lab. Since it could be too much if we applied many kinds of methods at the same time, which may decrease the resolution of images.

### IV. Analyze on the experiment results

#### A. What did you explore during the training process?

Both of the networks converge very fast on the training dataset, especially the UNet network, they also achieve good performance on validation and testing dataset.

```

100%|
[Epoch 100] loss = 0.017158072, average dice score = 0.986313522
100%|
Validation dice score = 0.9063050150871277

```

(a)

```

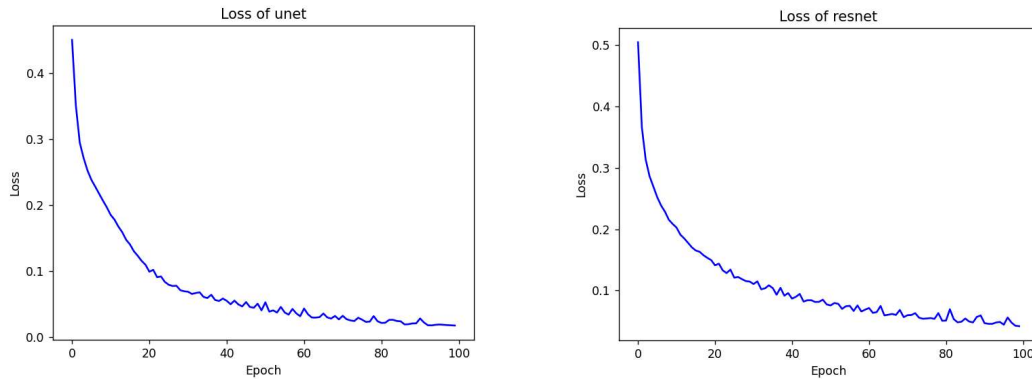
100%|
[Epoch 100] loss = 0.042193040, average dice score = 0.966403365
100%|
Validation dice score = 0.9000144600868225

```

(b)

**Fig. 8 The dice score of UNet (a) and ResNet34\_UNet (b) in the last epoch. (Both networks trained for 100 epochs with batch size 2 and learning rate 0.00001)**

From the following figures, we can see that though both of the networks achieved a very low loss, the UNet only needs 20 epochs to converge, while the ResNet34\_UNet needs 30 epochs or more. I think that may be because of the size of the network. Larger or deeper networks usually need more time to converge than smaller networks, since they have more weights to train, which may cause gradient vanishing or exploding problems.



**Fig. 9 Loss curve of UNet (left) and ResNet34\_UNet (right). The loss of UNet can be converged to about 0.01, while ResNet34\_UNet converges to 0.04**

Moreover, I observed that results of UNet (Fig.10-a) are sharper than the ResNet34\_UNet (Fig.10-b). In the below figure, we can see that although two networks predicted well on these images, the prediction of ResNet34\_UNet seems to have some blur.

Original Image



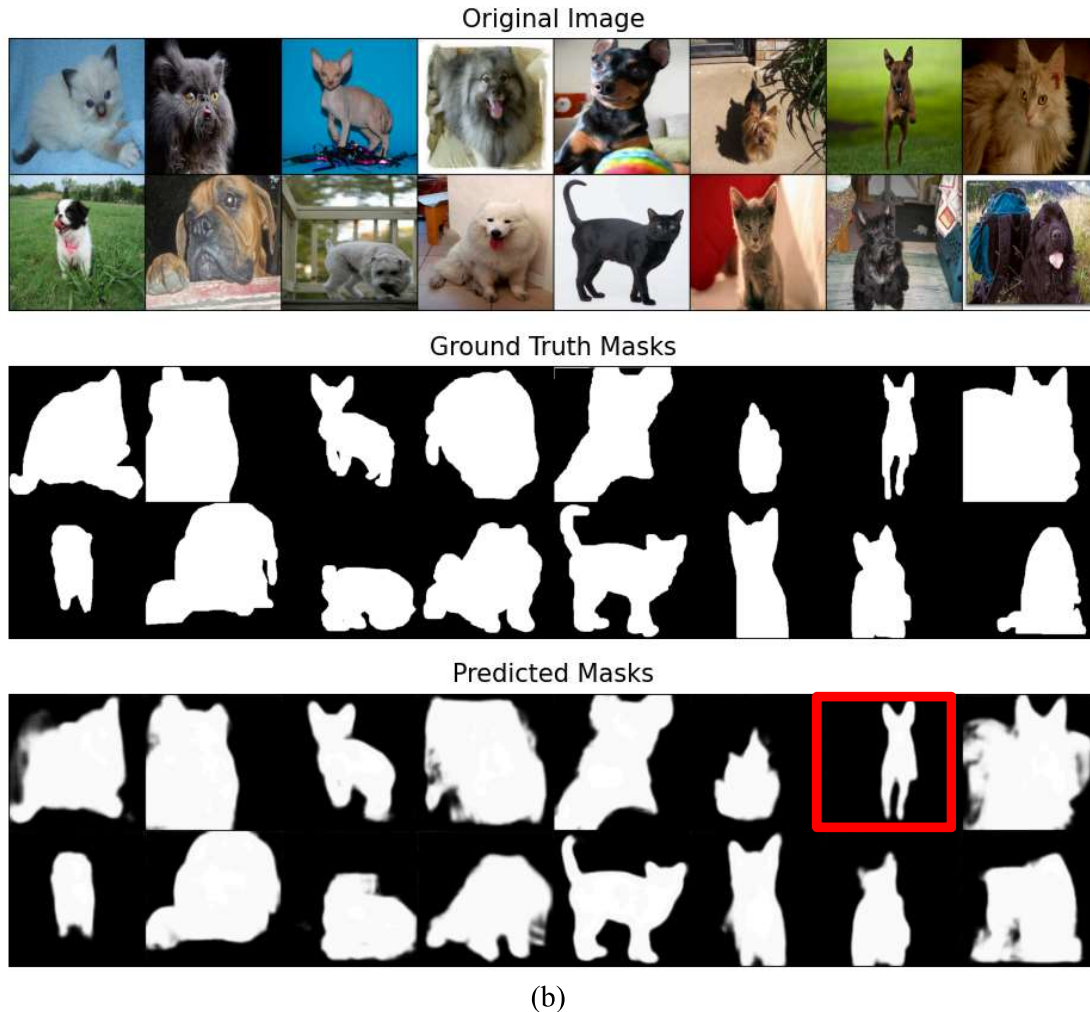
Ground Truth Masks



Predicted Masks



(a)



**Fig. 10 Predicted Masks of UNet (a) and ResNet34\_UNet (b).** We can apparently observe that UNet predicts a sharper mask.

## B. Found any characteristics of the data?

Objects in the images have large variations in pose, lighting, and scale, this characteristic allows networks to learn a more general feature to the foreground class, thus can achieve a similar score in the training and inferencing section.

## V. Execution command

For an anaconda environment, if you are at the root of the folder, please use the following command to train the model.



### A. The command and parameters for the training process

Model	Command
UNet	<code>python src/train.py -m unet</code>
ResNet34-UNet	<code>python src/train.py -m resnet</code>

Using the command above will train 50 epochs with the batch size of 2.

Parameters	Description	Default Value
-e	number of training epochs	50
-b	batch size	2
-lr	learning rate	1e-5
-m	model to train (options: unet/resnet) <b>(required)</b>	resnet
-d	device (option: cuda/cpu)	cuda
-p	save path of model checkpoints	./saved_models/
-r	name of model to keep on training (ex. unet.pth)	empty string

### B. The command and parameters for the inference process

Model	Command and parameters
UNet	<code>python src/inference.py -m DL_Lab3_UNet_313551073_顏琦恩.pth -u unet -b 16</code>
ResNet34-UNet	<code>python src/inference.py -m DL_Lab3_ResNet34_UNet_313551073_顏琦恩.pth -u resnet -b 16</code>

Using the command above will inference with the batch size of 16.

Parameters	Description	Default Value
-m	path to the stored model weight (ex. resnet.pth)	DL_Lab3_UNet_313551073_顏琦恩.pth
-b	batch size	2
-d	device (option: cuda/cpu)	cuda
-u	model to inference (options: unet/resnet) <b>(required)</b>	resnet

-p	plot during inferencing or not (option: y/n), it will show the batch which has the highest score before the batch inferencing	n
----	---	---

## VI. Discussion

### A. What architecture may bring better results?

In my experiment, the UNet architecture has better results than the ResNet34-UNet, but I think if there is more training data, or the testing data is more complicated, the ResNet34-UNet would get a better result, since the network is more deeper to learn a more complicated feature.

### B. What are the potential research topics in this task?

I think we can investigate the impact and effectiveness of attention mechanisms, such as CBAM, on image segmentation. Additionally, we can explore whether adding attention mechanisms into UNet can help improve accuracy.