

Part A: Pooling and Upsampling

1. The code for this question is shown below.

```

1  class PoolUpsampleNet(nn.Module):
2      def __init__(self, kernel, num_filters, num_colours, num_in_channels):
3          super().__init__()
4
5          # Useful parameters
6          padding = kernel // 2
7
8          ##### YOUR CODE GOES HERE #####
9          # construct cnn
10         group_one = nn.Sequential(
11             nn.Conv2d(in_channels=num_in_channels, out_channels=num_filters,
12                 kernel_size=kernel, padding=padding),
13             nn.MaxPool2d(kernel_size=2),
14             nn.BatchNorm2d(num_features=num_filters),
15             nn.ReLU()
16         )
17         group_two = nn.Sequential(
18             nn.Conv2d(in_channels=num_filters, out_channels=2*num_filters,
19                 kernel_size=kernel, padding=padding),
20             nn.MaxPool2d(kernel_size=2),
21             nn.BatchNorm2d(num_features=2*num_filters),
22             nn.ReLU()
23         )
24         group_three = nn.Sequential(
25             nn.Conv2d(in_channels=2*num_filters, out_channels=num_filters,
26                 kernel_size=kernel, padding=padding),
27             nn.Upsample(scale_factor=2),
28             nn.BatchNorm2d(num_features=num_filters),
29             nn.ReLU()
30         )
31         group_four = nn.Sequential(
32             nn.Conv2d(in_channels=num_filters, out_channels=num_colours,
33                 kernel_size=kernel, padding=padding),
34             nn.Upsample(scale_factor=2),
35             nn.BatchNorm2d(num_features=num_colours),
36             nn.ReLU()
37         )
38         group_five = nn.Sequential(
39             nn.Conv2d(in_channels=num_colours, out_channels=num_colours,
40                 kernel_size=kernel, padding=padding)
41         )
42         self.modules_list = nn.ModuleList([group_one, group_two, group_three,
43             group_four, group_five])
44
45         #####
46
47         def forward(self, x):
48             ##### YOUR CODE GOES HERE #####
49             for m in self.modules_list:
50                 x = m(x)
51             return x
52         #####

```

2. After 25 epochs using a batch size of 100 and a kernel size of 3, the best validation accuracy is 41.0%. Generally, these results do not look good although it uses different colors to color different objects. For example, according to Figure 1, this model performs badly for the background of the picture like the sky and the ground as it always presents gray sky and green ground.

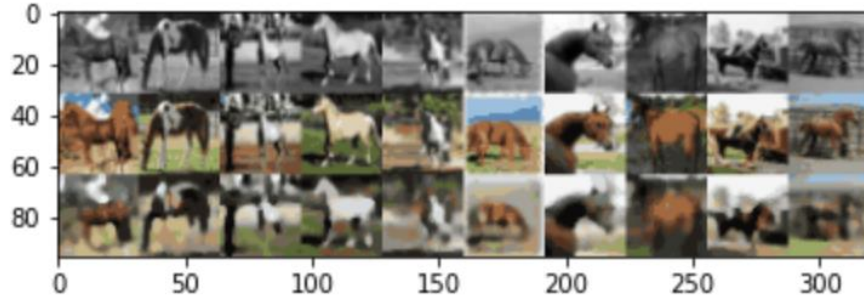


Figure 1: The trained result for some images.

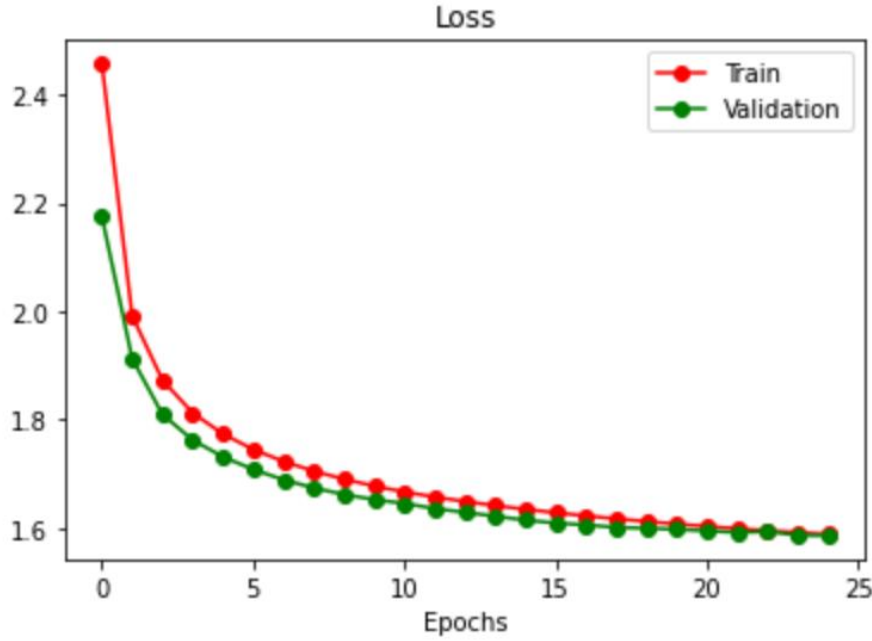


Figure 2: The training and validation loss.

- For the calculation of the convolution layer, I follow the instruction from the lecture four slide. The number of weights of the max pooling layer, ReLU layer, and up sampling layer are zero since there is no learning parameter. For this model, since the kernel size of the max pooling layer is 2×2 , the number of connections is basically the product of the output size and the kernel. Similarly, for the up sampling layer, the number of connections would be its output size because it just scales the tensor by 2. For the batch normalization layer, since the mean and variance of the normalization depends on the feature number, the number of connections is the product of the output size and the size of the height and width. To make the it clear, let's write it in a table.

Thus, the total number of weights is

$$\begin{aligned}
 & 9NIC * NF + NF + 2NF + 18NF^2 + 2NF + 4NF + 18NF^2 + NF + 2NF + 9NF * NC + NC + 2NC + 9NC^2 \\
 & = 36NF^2 + 9NC^2 + 9NF * NC + 9NIC * NF + 12NF + 4NC
 \end{aligned}$$

The number of outputs is $32 * 32 * NC = 1024NC$.

Layer	Input	Output	# weights	# outputs	# connections
Image	-	$32 * 32 * 1$	-	-	-
Conv2d	$32 * 32 * NIC$	$32 * 32 * NF$	$3^2 * NIC * NF + NF$	$32 * 32 * NF$	$32 * 32 * 3^2 * NF * NIC$
MaxPool2d	$32 * 32 * NF$	$16 * 16 * NF$	0	$16 * 16 * NF$	$16 * 16 * 2^2 * NF$
BatchNorm2d	$16 * 16 * NF$	$16 * 16 * NF$	$2 * NF$	$16 * 16 * NF$	$16 * 16 * NF * 16^2$
ReLU	$16 * 16 * NF$	$16 * 16 * NF$	0	$16 * 16 * NF$	-
Conv2d	$16 * 16 * NF$	$16 * 16 * 2NF$	$3^2 * NF * 2NF + 2NF$	$16 * 16 * 2NF$	$16 * 16 * 3^2 * NF * 2NF$
MaxPool2d	$16 * 16 * 2NF$	$8 * 8 * 2NF$	0	$8 * 8 * 2NF$	$8 * 8 * 2^2 * 2NF$
BatchNorm2d	$8 * 8 * 2NF$	$8 * 8 * 2NF$	$4 * NF$	$8 * 8 * 2NF$	$8 * 8 * 2NF * 8^2$
ReLU	$8 * 8 * 2NF$	$8 * 8 * 2NF$	0	$8 * 8 * 2NF$	-
Conv2d	$8 * 8 * 2NF$	$8 * 8 * NF$	$3^2 * NF * 2NF + NF$	$8 * 8 * NF$	$8 * 8 * 3^2 * NF * 2NF$
Upsample	$8 * 8 * NF$	$16 * 16 * NF$	0	$16 * 16 * NF$	$8 * 8 * NF * 2^2$
BatchNorm2d	$16 * 16 * NF$	$16 * 16 * NF$	$2 * NF$	$16 * 16 * NF$	$16 * 16 * NF * 16^2$
ReLU	$16 * 16 * NF$	$16 * 16 * NF$	0	$16 * 16 * NF$	-
Conv2d	$16 * 16 * NF$	$16 * 16 * NC$	$3^2 * NF * NC + NC$	$16 * 16 * NC$	$16 * 16 * 3^2 * NF * NC$
Upsample	$16 * 16 * NC$	$32 * 32 * NC$	0	$32 * 32 * NC$	$16 * 16 * NF * 2^2$
BatchNorm2d	$32 * 32 * NC$	$32 * 32 * NC$	$2 * NC$	$32 * 32 * NC$	$32 * 32 * NC * 32^2$
ReLU	$32 * 32 * NC$	$32 * 32 * NC$	0	$32 * 32 * NC$	-
Conv2d	$32 * 32 * NC$	$32 * 32 * NC$	$3^2 * NC * NC + NC$	$32 * 32 * NC$	$32 * 32 * 3^2 * NC * NC$

The total number of connections is

$$\begin{aligned}
& 9 * 1024NF * NIC + 4 * 256NF + 256 * 256NF + 18 * 256NF^2 + 4 * 128NF + 2 * 64 * 64NF \\
& + 18 * 64NF^2 + 4 * 64NF + 256 * 256NF + 9 * 256NF * NC + 4 * 256NF + 1024 * 1024NC + 9 * 1024NC^2 \\
& = 9216NF * NIC + 9216NC^2 + 5760NF^2 + 2304NF * NC + 142080NF + 1048576NC
\end{aligned}$$

Similarly, for the doubled input dimension, we can get the following table.

Layer	Input	Output	# weights	# outputs	# connections
Image	-	$64 * 64 * 1$	-	-	-
Conv2d	$64 * 64 * NIC$	$64 * 64 * NF$	$3^2 * NIC * NF + NF$	$64 * 64 * NF$	$64 * 64 * 3^2 * NF * NIC$
MaxPool2d	$64 * 64 * NF$	$32 * 32 * NF$	0	$32 * 32 * NF$	$32 * 32 * 2^2 * NF$
BatchNorm2d	$32 * 32 * NF$	$32 * 32 * NF$	$2 * NF$	$32 * 32 * NF$	$32 * 32 * NF * 32^2$
ReLU	$32 * 32 * NF$	$32 * 32 * NF$	0	$32 * 32 * NF$	-
Conv2d	$32 * 32 * NF$	$32 * 32 * 2NF$	$3^2 * NF * 2NF + 2NF$	$32 * 32 * 2NF$	$32 * 32 * 3^2 * NF * 2NF$
MaxPool2d	$32 * 32 * 2NF$	$16 * 16 * 2NF$	0	$16 * 16 * 2NF$	$16 * 16 * 2^2 * 2NF$
BatchNorm2d	$16 * 16 * 2NF$	$16 * 16 * 2NF$	$4 * NF$	$16 * 16 * 2NF$	$16 * 16 * 2NF * 16^2$
ReLU	$16 * 16 * 2NF$	$16 * 16 * 2NF$	0	$16 * 16 * 2NF$	-
Conv2d	$16 * 16 * 2NF$	$16 * 16 * NF$	$3^2 * NF * 2NF + NF$	$16 * 16 * NF$	$16 * 16 * 3^2 * NF * 2NF$
Upsample	$16 * 16 * NF$	$32 * 32 * NF$	0	$32 * 32 * NF$	$16 * 16 * NF * 2^2$
BatchNorm2d	$32 * 32 * NF$	$32 * 32 * NF$	$2 * NF$	$32 * 32 * NF$	$32 * 32 * NF * 32^2$
ReLU	$32 * 32 * NF$	$32 * 32 * NF$	0	$32 * 32 * NF$	-
Conv2d	$32 * 32 * NF$	$32 * 32 * NC$	$3^2 * NF * NC + NC$	$32 * 32 * NC$	$32 * 32 * 3^2 * NF * NC$
Upsample	$32 * 32 * NC$	$64 * 64 * NC$	0	$64 * 64 * NC$	$32 * 32 * NF * 2^2$
BatchNorm2d	$64 * 64 * NC$	$64 * 64 * NC$	$2 * NC$	$64 * 64 * NC$	$64 * 64 * NC * 64^2$
ReLU	$64 * 64 * NC$	$64 * 64 * NC$	0	$64 * 64 * NC$	-
Conv2d	$64 * 64 * NC$	$64 * 64 * NC$	$3^2 * NC * NC + NC$	$64 * 64 * NC$	$64 * 64 * 3^2 * NC * NC$

Thus, the total number of weights is

$$\begin{aligned}
& 9NIC * NF + NF + 2NF + 18NF^2 + 2NF + 4NF + 18NF^2 + NF + 2NF + 9NF * NC + NC + 2NC + 9NC^2 \\
& = 36NF^2 + 9NC^2 + 9NF * NC + 9NIC * NF + 12NF + 4NC
\end{aligned}$$

The number of outputs is $64 * 64 * NC = 4096NC$.

The total number of connections is

$$\begin{aligned}
& 9 * 4096NF * NIC + 4 * 1024NF + 1024 * 1024NF + 18 * 1024NF^2 + 4 * 512NF + 2 * 256 * 256NF \\
& + 18 * 256NF^2 + 4 * 256NF + 1024 * 1024NF + 9 * 1024NF * NC + 4 * 1024NF + 4096 * 4096NC + 9 * 4096NC^2 \\
& = 36864NF * NIC + 36864NC^2 + 23040NF^2 + 9216NF * NC + 142080NF + 2239488NC
\end{aligned}$$

Part B: Strided and Transposed Convolutions

1. The code for this question is shown below.

```

1  class ConvTransposeNet(nn.Module):
2      def __init__(self, kernel, num_filters, num_colours, num_in_channels):
3          super().__init__()
4
5          # Useful parameters
6          stride = 2
7          padding = kernel // 2
8          output_padding = 1
9
10         ##### YOUR CODE GOES HERE #####
11         group_one = nn.Sequential(
12             nn.Conv2d(in_channels=num_in_channels, out_channels=num_filters,
13                 kernel_size=kernel, stride=2, padding=1),
14             nn.BatchNorm2d(num_features=num_filters),
15             nn.ReLU()
16         )
17         group_two = nn.Sequential(
18             nn.Conv2d(in_channels=num_filters, out_channels=2*num_filters,
19                 kernel_size=kernel, stride=2, padding=1),
20             nn.BatchNorm2d(num_features=2*num_filters),
21             nn.ReLU()
22         )
23         group_three = nn.Sequential(
24             nn.ConvTranspose2d(in_channels=2*num_filters, out_channels=num_filters,
25                 kernel_size=kernel, stride=2, padding=1, output_padding=1),
26             nn.BatchNorm2d(num_features=num_filters),
27             nn.ReLU()
28         )
29         group_four = nn.Sequential(
30             nn.ConvTranspose2d(in_channels=num_filters, out_channels=num_colours,
31                 kernel_size=kernel, stride=2, padding=1, output_padding=1),
32             nn.BatchNorm2d(num_features=num_colours),
33             nn.ReLU()
34         )
35         group_five = nn.Sequential(
36             nn.Conv2d(in_channels=num_colours, out_channels=num_colours,
37                 kernel_size=kernel, padding=padding)
38         )
39         self.modules_list = nn.ModuleList([group_one, group_two, group_three,
40             group_four, group_five])
41
42         ##### YOUR CODE GOES HERE #####
43         def forward(self, x):
44             ##### YOUR CODE GOES HERE #####
45             for m in self.modules_list:
46                 x = m(x)
47             return x
48         ##### YOUR CODE GOES HERE #####

```

2. After 25 epochs using a batch size of 100 and a kernel size of 3, the best validation accuracy is 55.7%.
3. The result is better than Part A. According to Figure 3, this model performs better for the background of the picture like the sky, but it still always presents green ground. According to



Figure 3: The trained result for some images.

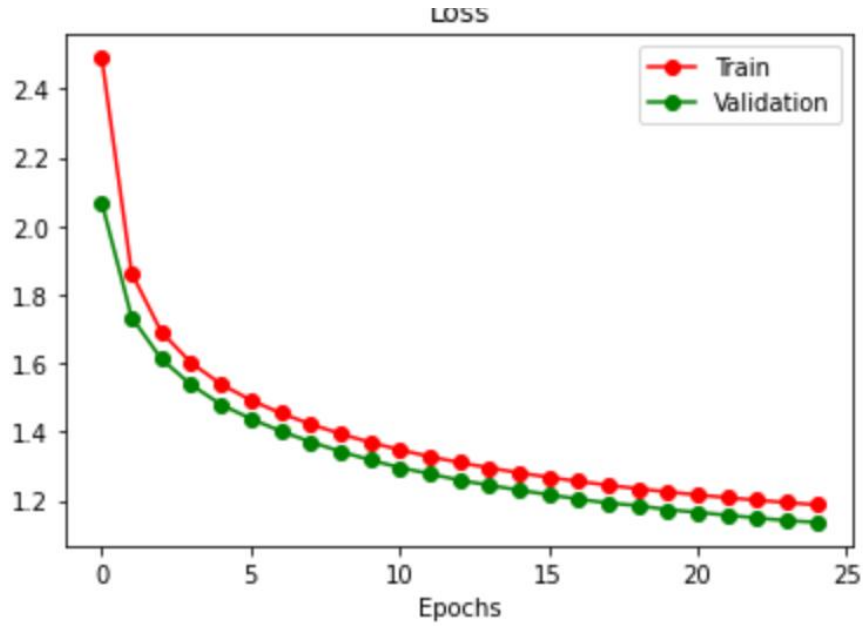


Figure 4: The training and validation loss.

Figure 2 and 4, we can see that the validation loss is always smaller than the training loss for the *ConvTransposeNet* model. Compared to the *PoolUpsampleNet* model, the validation loss in the *ConvTransposeNet* model is lower.

4. Let's consider the padding parameter for the first two convolution layers first. According to the previous setting, we know $padding = 1$ and $stride = 2$ for the case of $kernel_size = 3$. For the dilation parameter, we use the default setting which is 1. Thus, define a be the input size, b be the output size, p be the padding size, and k be the kernel size. We can get

$$\begin{aligned}
 b &= \frac{a + 2p - 1 * (k - 1) - 1}{2} + 1 \\
 &= \frac{a + 2 * 1 - (3 - 1) - 1}{2} + 1 \\
 &= \frac{a + 1}{2}
 \end{aligned}$$

To maintain the output size, it must follow

$$\begin{aligned} 2p - (k - 1) + 1 &= 1 \\ \implies 2p - k &= -1 \\ \implies p &= \frac{k - 1}{2} \end{aligned}$$

Thus, when $k = 4$, $p = 1.5$. However, 1.5 is invalid for the number of padding. Thus, we should either set $p = 1$ or $p = 2$. When $k = 5$, $p = 2$.

For the padding parameters of the *ConvTranspose2d* layer, we know *padding* = 1, *stride* = 2, and *output_padding* = 1 for the case of *kernel_size* = 3. Similarly, for the dilation parameter, we use the default setting which is 1. Define *o* be the size of output padding. We can get

$$\begin{aligned} b &= s(a - 1) - 2p + 1 * (k - 1) + o + 1 \\ &= 2a - 2 - 2 + 2 + 1 + 1 \\ &= 2a \end{aligned}$$

To maintain the output size, it must follow

$$\begin{aligned} -2p + k + o - 2 &= 0 \\ \implies 2p - o &= k - 2 \end{aligned}$$

Thus, there are many solutions. When $k = 4$, $2p - o = 2$. For example, we can set *padding* = 1 and *output_padding* = 0. When $k = 5$, $2p - o = 3$. Setting *padding* = 2 and *output_padding* = 1 works.

5. I set the batch size be 32, 64, 128, 256, and 512 with 25 epochs. The following table shows the corresponding best validation loss.

Batch size	Best validation loss
32	57.4%
64	55.2%
128	53.1%
256	50.6%
512	45.9%

According to the result, we can see that the model performs worse as we increase the batch size. When the batch size is 32, it achieved the highest accuracy as well as the best final image output comparing with others. As the batch size becomes larger, the difference between the training loss and the validation loss becomes smaller. Thus, we can conclude that the model prefers small value of the batch size.

Part C: Skip Connections

1. The code for this question is shown below.

```

1 class UNet(nn.Module):
2     def __init__(self, kernel, num_filters, num_colours, num_in_channels):
3         super().__init__()
4
5         # Useful parameters
6         stride = 2
7         padding = kernel // 2
8         output_padding = 1
9
10        ##### YOUR CODE GOES HERE #####
11        self.group_one = nn.Sequential(
12            nn.Conv2d(in_channels=num_in_channels, out_channels=num_filters,
13                    kernel_size=kernel, stride=stride, padding=padding),

```

```

13         nn.BatchNorm2d(num_features=num_filters),
14         nn.ReLU()
15     )
16     self.group_two = nn.Sequential(
17         nn.Conv2d(in_channels=num_filters, out_channels=2*num_filters,
18 kernel_size=kernel, stride=stride, padding=padding),
19         nn.BatchNorm2d(num_features=2*num_filters),
20         nn.ReLU()
21     )
22     self.group_three = nn.Sequential(
23         nn.ConvTranspose2d(in_channels=2*num_filters, out_channels=num_filters
24 , kernel_size=kernel, stride=stride, padding=padding, output_padding=
25 output_padding),
26         nn.BatchNorm2d(num_features=num_filters),
27         nn.ReLU()
28     )
29     self.group_four = nn.Sequential(
30         nn.ConvTranspose2d(in_channels=2*num_filters, out_channels=num_colours
31 , kernel_size=kernel, stride=2, padding=padding, output_padding=output_padding
32 ),
33         nn.BatchNorm2d(num_features=num_colours),
34         nn.ReLU()
35     )
36     self.group_five = nn.Sequential(
37         nn.Conv2d(in_channels=num_in_channels+num_colours, out_channels=
38 num_colours, kernel_size=kernel, padding=padding)
39     )
40     #####
41
42 def forward(self, x):
43     ##### YOUR CODE GOES HERE #####
44     group_one_out = self.group_one(x)
45     group_two_out = self.group_two(group_one_out)
46     group_three_out = self.group_three(group_two_out)
47     group_four_in = torch.cat((group_one_out, group_three_out), 1)
48     group_four_out = self.group_four(group_four_in)
49     group_five_in = torch.cat((x, group_four_out), 1)
50     return self.group_five(group_five_in)
51     #####

```

2. After 25 epochs using a batch size of 100 and a kernel size of 3, the accuracy is 58.2%.

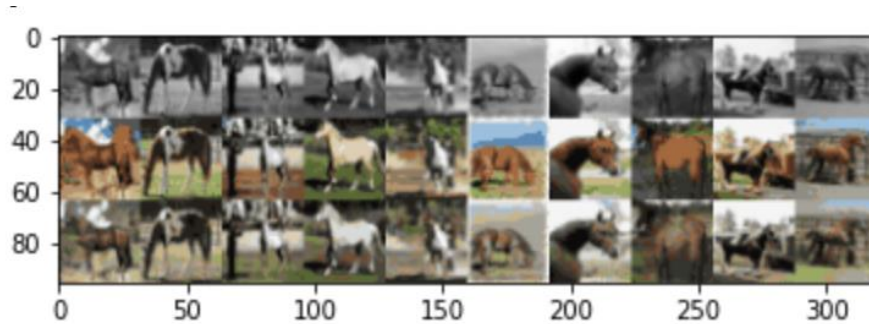


Figure 5: The trained result for some images.

3. The result is better than both Part A and Part B. Thus, we can conclude that the skip connections improve the accuracy and the output qualitatively. According to Figure 4, this model performs better for the background of the picture like the sky, and it performs better about the ground as it uses different colors to paint ground. According to Figure 2, 4 and 6, we can see that the validation loss is always smaller than the training loss for the *UNet* model. Compared to the *PoolUpsampleNet* and *ConvTransposeNet* model, the validation loss in the *UNet* model is

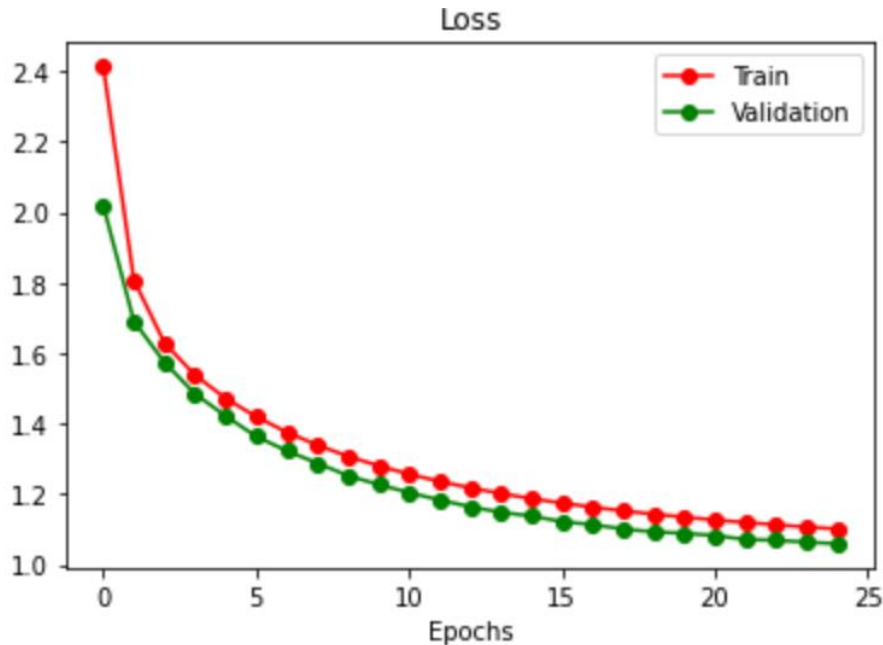


Figure 6: The training and validation loss.

lower. I think one of the reasons is that skip connections require more trainable parameters since we concatenate two layers in a pattern, which can learn more things from the input and perform better. The other reason might be concatenating two layers makes up for the possible missing information from the pooling or sampling process.

Part D.1: Fine-tune Semantic Segmentation Model with Cross Entropy Loss

1. The code for this question is shown below.

```

1 def train(args, model):
2
3     # Set the maximum number of threads to prevent crash in Teaching Labs
4     torch.set_num_threads(5)
5     # Numpy random seed
6     np.random.seed(args.seed)
7
8     # Save directory
9     # Create the outputs folder if not created already
10    save_dir = "outputs/" + args.experiment_name
11    if not os.path.exists(save_dir):
12        os.makedirs(save_dir)
13
14    learned_parameters = []
15    # We only learn the last layer and freeze all the other weights
16    ##### Code goes here #####
17    # Around 3 lines of code
18    # Hint:
19    # - use a for loop to loop over all model.named_parameters()
20    # - append the parameters (both weights and biases) of the last layer (prefix:
21    #   classifier.4) to the learned_parameters list
22    for name, param in model.named_parameters():
23        if name.startswith("classifier.4"):
24            learned_parameters.append(param)
25    #####

```

2. The code for this question is shown below.


```

1 # Truncate the last layer and replace it with the new one.
2 # To avoid 'CUDA out of memory' error, you might find it useful (sometimes
   required)
3 # to set the 'requires_grad' = False for some layers
4 ##### YOUR CODE GOES HERE #####
5 # Around 4 lines of code
6 # Hint:
7 # - replace the classifier.4 layer with the new Conv2d layer (1 line)
8 # - no need to consider the aux_classifier module (just treat it as don't care)
9 # - freeze the gradient of other layers (3 lines)
10 model.classifier._modules['4'] = nn.Conv2d(256, 2, kernel_size=(1, 1), stride=(1,
   1))
11 for name, param in model.named_parameters():
12     if not name.startswith("classifier.4"):
13         param.requires_grad = False
14 #####
15

```

According to the output, the best validation mIoU is 0.3373.

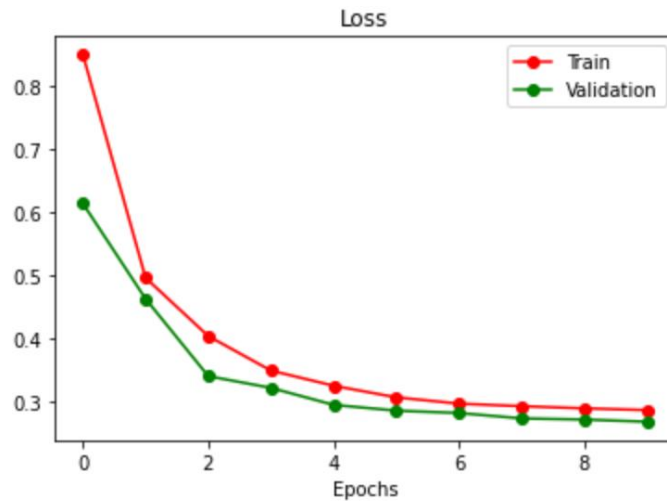


Figure 7: The training and validation loss.

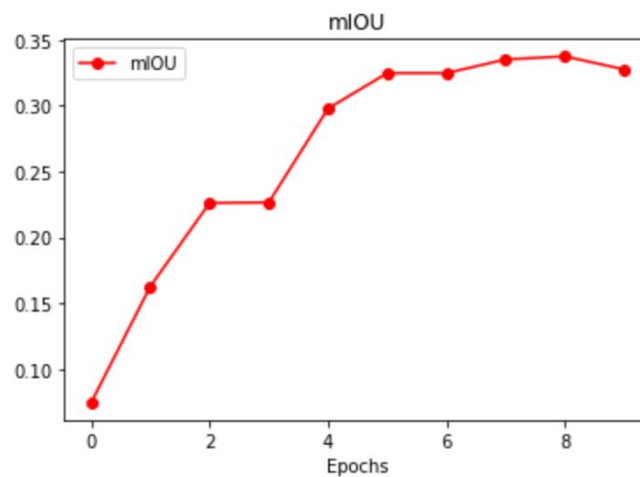


Figure 8: The mIoU.

3. Visualization

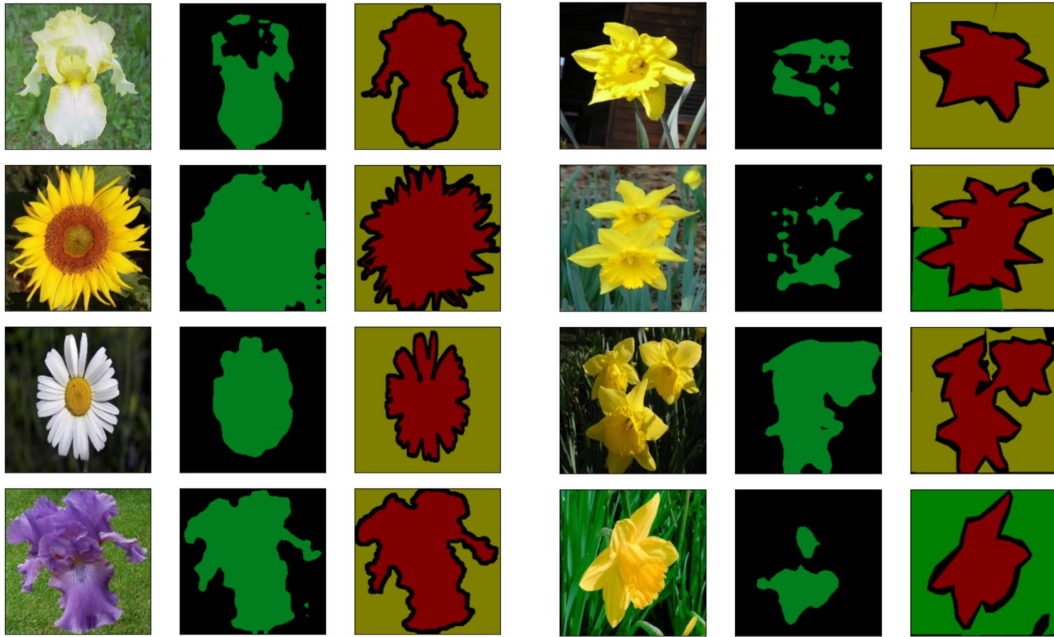


Figure 9: The predictions visualization.

Part D.2: Fine-tune Semantic Segmentation Model with IoU Loss

1. The code for this question is shown below.

```

1 def compute_iou_loss(pred, gt, SMOOTH=1e-6):
2     # Compute the IoU between the pred and the gt (ground truth)
3     ##### YOUR CODE GOES HERE #####
4     # Around 5 lines of code
5     # Hint:
6     # - apply softmax on pred along the channel dimension (dim=1)
7     # - only have to compute IoU between gt and the foreground channel of pred
8     # - no need to consider IoU for the background channel of pred
9     # - extract foreground from the softmaxed pred (e.g., softmaxed_pred[:, 1, :, :])
10    # - compute intersection between foreground and gt
11    # - compute union between foreground and gt
12    # - compute loss using the computed intersection and union
13    softmaxed_pred = F.softmax(pred, dim=1)
14    new_pred = softmaxed_pred[:, 1, :, :]
15    i = torch.mul(new_pred, gt)
16    u = new_pred + gt - i
17    loss = 1 - torch.sum(i)/torch.sum(u)
18    #####
19
20    return loss
21

```

According to the output, the best validation mIoU is 0.3014. It performs slightly worse than the previous model although we use stronger loss function. However, according to Figure 10, we can see that the corresponding training loss and validation loss are about 0.4, which seems not to be converged and could be improved. Thus, if we tune some hyperparameters like the number of epochs, we might get better result.

2. Visualization

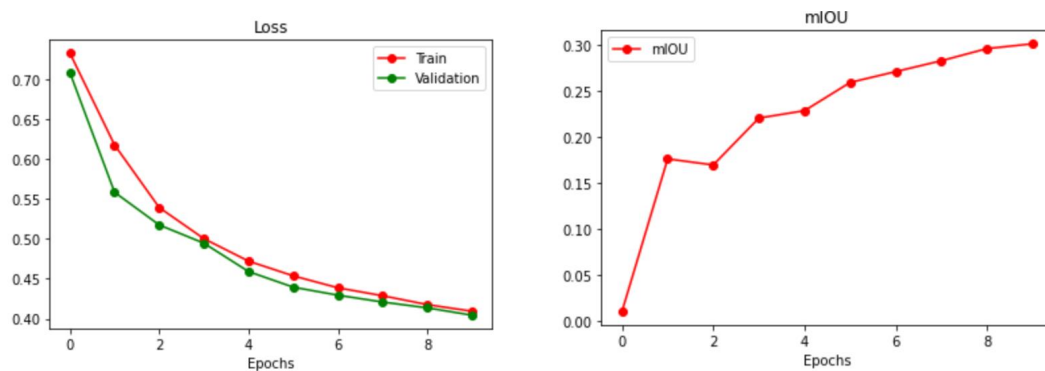


Figure 10: The loss and mIoU.

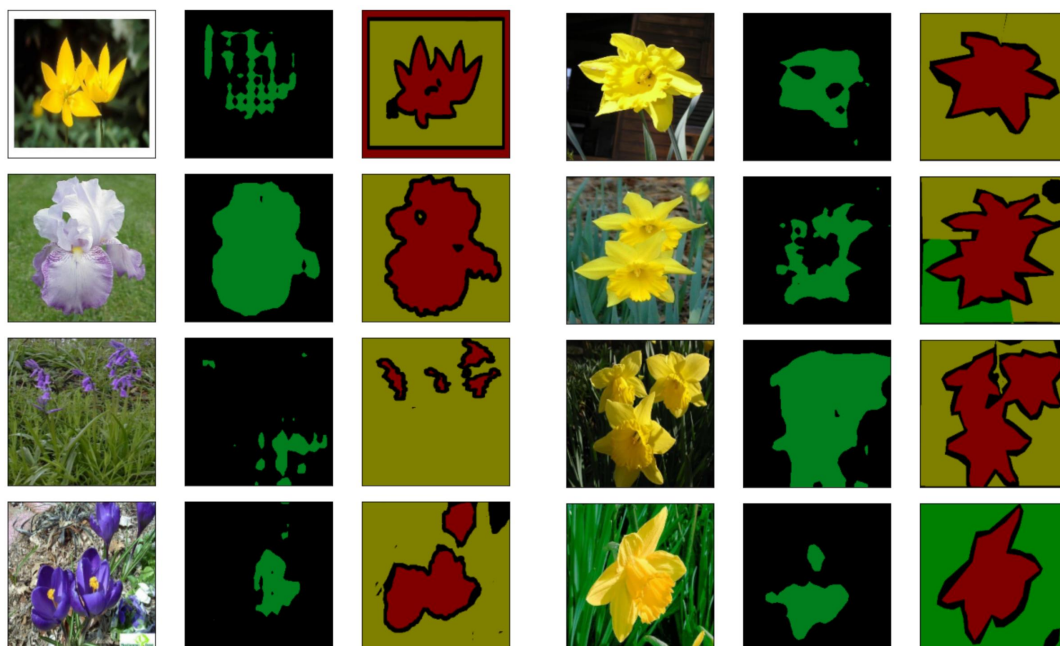


Figure 11: The predictions visualization.