# Deep Learning for Computer Vision
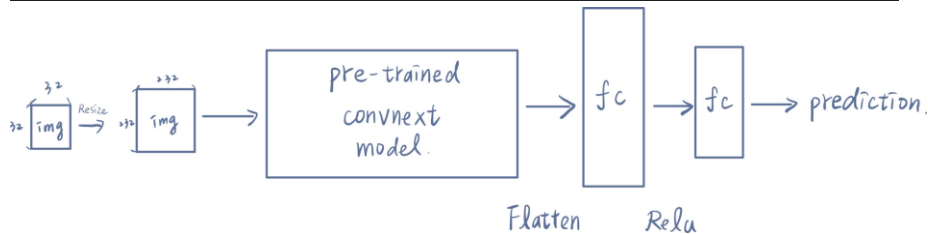
電機所碩一 謝宗翰 R12921A10

- **Problem 1: Image Classification** (25%)

1. *Draw the network architecture of method A or B*

```
1  model = convnext_base(weights=ConvNeXt_Base_Weights.DEFAULT)
2
3  model.classifier = nn.Sequential(
4      nn.Flatten(), nn.Linear(1024, 500), nn.ReLU(), nn.Linear(500, 50)
5  )
```



2. *Report accuracy of your models (both A, B) on the validation set.*

| Accuracy | |
|---|---|
| A | B |
| 30.7% | 90.2% |

3. *Report your implementation details of model A.*

```
1  class CNN(nn.Module):
2      def __init__(self):
3          super(CNN, self).__init__()  # 可調用nn.Moudule的函數
4
5          # 第一個卷積層
6          self.conv1 = nn.Conv2d(in_channels=3, out_channels=32, kernel_size=3, padding=1)
7          self.relu1 = nn.ReLU()
8          self.pool1 = nn.MaxPool2d(kernel_size=2, stride=2)
9
10         # 第二個卷積層
11         self.conv2 = nn.Conv2d(
12             in_channels=32, out_channels=64, kernel_size=3, padding=1
13         )
14         self.relu2 = nn.ReLU()
15         self.pool2 = nn.MaxPool2d(kernel_size=2, stride=2)
16
17         # 全連接層
18         self.fc1 = nn.Linear(in_features=29696, out_features=128)
19         self.relu3 = nn.ReLU()
20         self.fc2 = nn.Linear(in_features=128, out_features=50)  # 50類
21
22         self.drop25 = nn.Dropout2d(0.25)
```

```
23
24      def forward(self, x):
25          x = self.conv1(x)
26          x = self.relu1(x)
27          x = self.pool1(x)
28
29          x = self.conv2(x)
30          x = self.relu2(x)
31          x = self.pool2(x)
32
33          x = x.view(x.size(0), -1)   # 展平特徵圖
34          x = self.fc1(x)
35          return x
```

| A Model | | | | |
|---|---|---|---|---|
| Learning Rate | Batch Size | Optimizer | Epochs | Loss Func |
| 0.0003 | 32 | Adam | 80 | Cross Entropy |
| Train Data augmentation | | | Valid Data augmentation | |
| No | | | No | |

| B Model | | | | |
|---|---|---|---|---|
| Learning Rate | Batch Size | Optimizer | Epochs | Loss Func |
| 0.0005 | 128 | Adam | 50 | Cross Entropy |
| Train Data augmentation | | | | |

```
1   basic_transform = [
2       # additional data argument
3       trns.RandomResizedCrop(224, scale=(0.08, 1.0), ratio=(3.0 / 4.0, 4.0 / 3.0)),
4       trns.RandomHorizontalFlip(),  # 水平翻轉
5       trns.ColorJitter(brightness=0.5, contrast=0.5, hue=0.5),
6   ]
7
8   transform_train = trns.Compose(
9       [
10          trns.Resize((232, 232), interpolation=trns.InterpolationMode.BICUBIC),
11          trns.CenterCrop((224, 224)),
12          trns.RandomChoice(basic_transform),
13          trns.ToTensor(),
14          trns.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
15      ]
16  )
```

Valid Data augmentation

```
1   transform_val = trns.Compose(
2       [
3           trns.Resize((232, 232), interpolation=trns.InterpolationMode.BICUBIC),
4           trns.CenterCrop((224, 224)),
5           trns.ToTensor(),
6           trns.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
7       ]
8   )
```

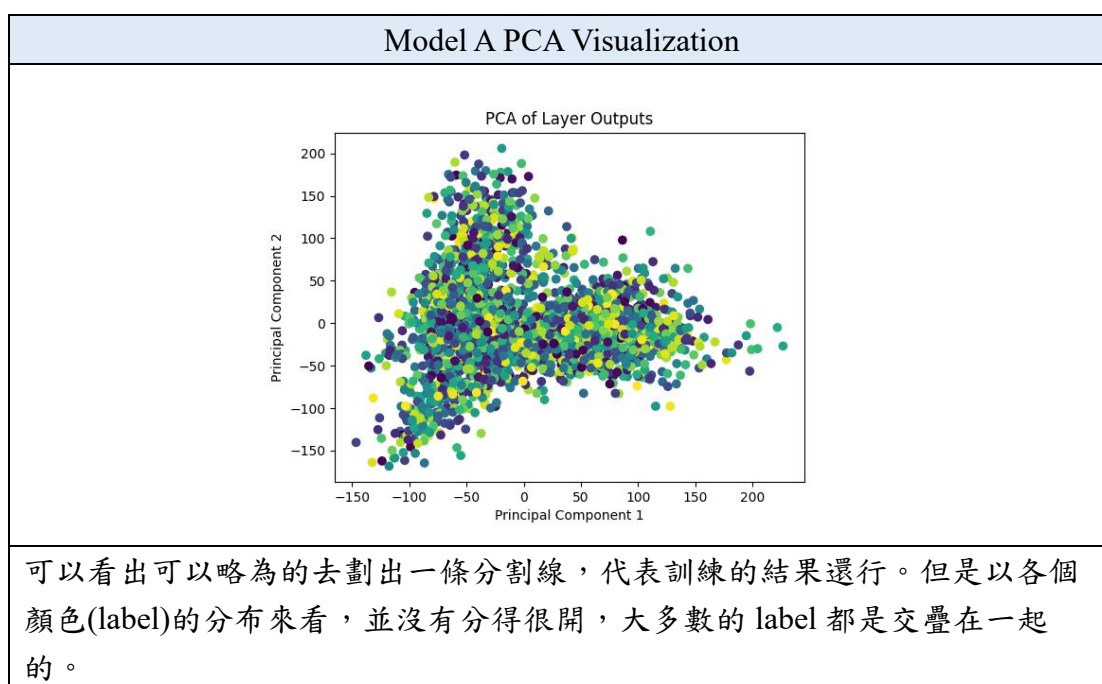4. *Report your alternative model or method in B, and describe its difference from model A.*

我的 B 用 torchvision model 的 ConvNext 進行 finetuning，把第 7 層前的

weight 都 freeze 住,只 train 倒數幾層捲積層及 Classifier。
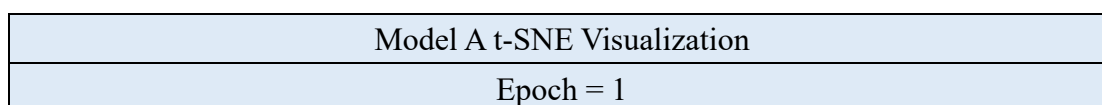
```
1  for i, (name, param) in enumerate(model.named_parameters()):
2      param.requires_grad = False
3      if "classifier" in name or "7" in name:
4          param.requires_grad = True
```
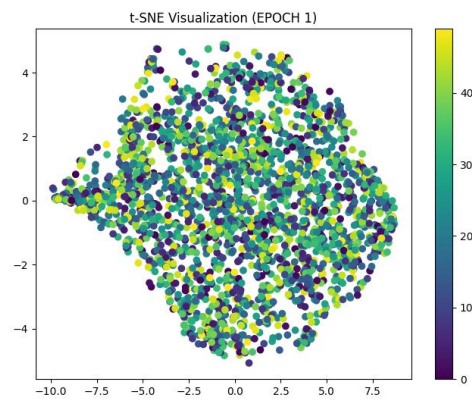
跟 A 比較不一樣的地方是,B 是用 pre-train 的 weight,著重在 Classifier 的訓練,讓 pre-train model 能對提供的 dataset 進行分類。B 模型比我自己手刻的 Net 深很多,效果比我的 A 好非常多。

5. *Visualize the learned visual representations of model A on the validation set by implementing PCA (Principal Component Analysis) on the output of the second last layer. Briefly explain your result of the PCA visualization.*

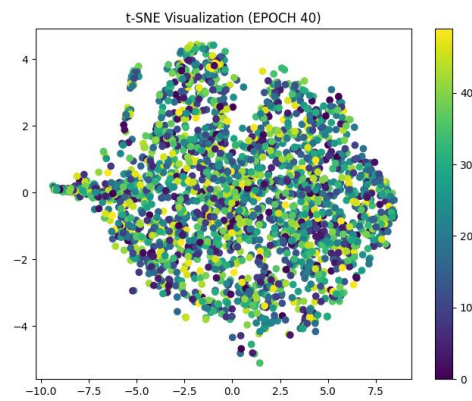| Model A PCA Visualization |
|---|
|  |
| 可以看出可以略為的去劃出一條分割線,代表訓練的結果還行。但是以各個顏色(label)的分布來看,並沒有分得很開,大多數的 label 都是交疊在一起的。 |

6. *Visualize the learned visual representation of model A, again on the output of the second last layer, but using t-SNE (t-distributed Stochastic Neighbor Embedding) instead. Depict your visualization from three different epochs including the first one and the last one. Briefly explain the above results.*
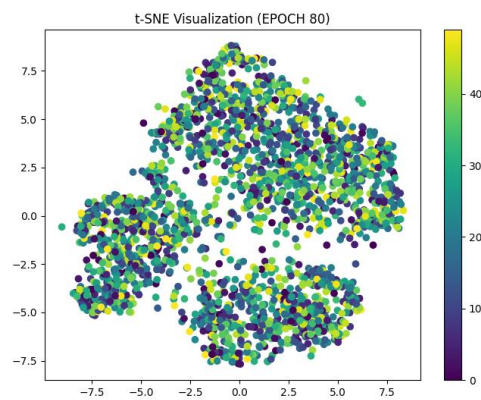
| Model A t-SNE Visualization |
|---|
| Epoch = 1 |

t-SNE Visualization (EPOCH 1)

| Epoch = 40 |
| --- |


t-SNE Visualization (EPOCH 40)

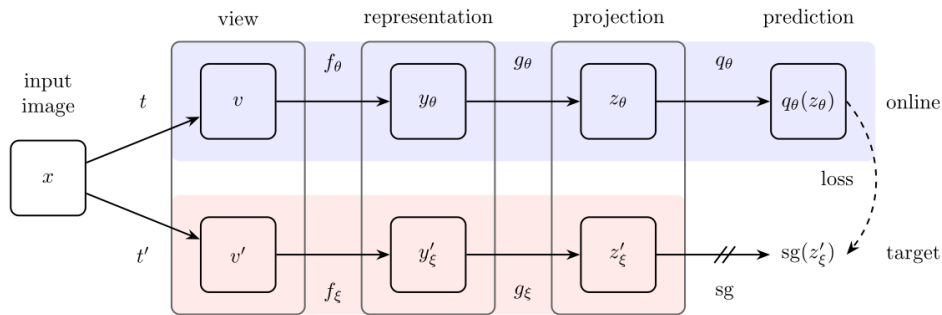| Epoch = 80 |
| --- |


t-SNE Visualization (EPOCH 80)

隨著 Epoch 越來越大,可以看出各點開始由中心開始往外擴,由訓練到 Epoch = 80 時非常明顯。但顏色跟顏色之間也都是交疊在一起,無明顯區別。

## ● Problem 2: Self-Supervised Pre-training for Image Classification

*1. Describe the implementation details of your SSL method for pre-training the ResNet50 backbone.*

我的 SSL 採用助教提供的方法(BYOL)



上圖為 BYOL 的架構,他有兩個網路,分別為 Target 及 Online

(1) 先對 Input x 做 data argmentation,得到 $t$ 及 $t'$

(2) 將 $t$ 輸入至 Online 網路中經$f_\theta$ 提取特徵,得到特徵向量$y_\theta$。同時將$t'$輸入到 target 中,經$f_\varepsilon$提取特徵,得到特徵向量$y'_\varepsilon$。

(3) $y_\theta$ 經過 MLP 網路$g_\theta$,得到$z_\theta$。而$y'_\varepsilon$ 經過 MLP 網路$g_\varepsilon$,得到$z'_\varepsilon$。

(4) $z_\theta$ 經過 MLP 網路$q_\theta$ ,得到$q_\theta(z_\theta)$,與$z'_\varepsilon$計算 loss

(5) *loss* 參照 BYOL 的論文

$$\mathcal{L}_{\theta,\xi} \triangleq \left\| \overline{q_\theta}(z_\theta) - \overline{z}'_\xi \right\|_2^2 = 2 - 2 \cdot \frac{\langle q_\theta(z_\theta), z'_\xi \rangle}{\left\| q_\theta(z_\theta) \right\|_2 \cdot \left\| z'_\xi \right\|_2}.$$
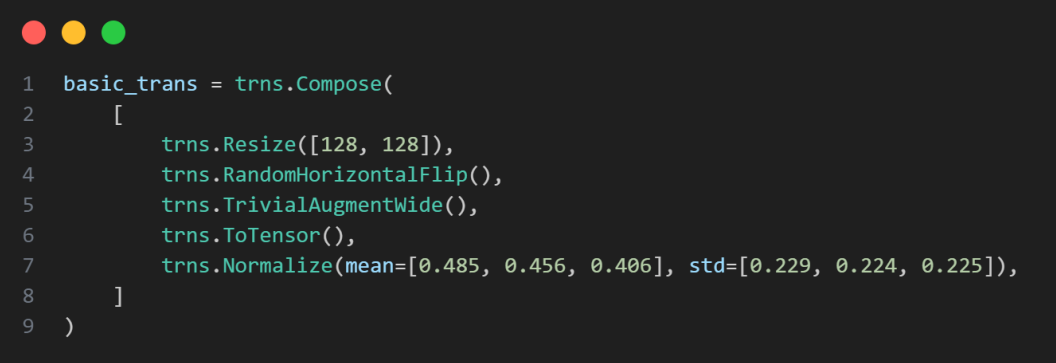
| SSL Model (BYOL) | | | | |
|---|---|---|---|---|
| Learning Rate | Batch Size | Optimizer | Epochs | Loss Func |
| 0.01 | 256 | Adam | 10 | Focal Loss |
| Train Data augmentation | | | | |

```
1  basic_transform = [
2      # additional data argument
3      trns.RandomResizedCrop(224, scale=(0.08, 1.0), ratio=(3.0 / 4.0, 4.0 / 3.0)),
4      trns.RandomHorizontalFlip(),  # 水平翻轉
5      trns.ColorJitter(brightness=0.5, contrast=0.5, hue=0.5),
6  ]
7
8  transform_train = trns.Compose(
9      [
10         trns.Resize((232, 232), interpolation=trns.InterpolationMode.BICUBIC),
11         trns.CenterCrop((224, 224)),
12         trns.RandomChoice(basic_transform),
13         trns.ToTensor(),
14         trns.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
15     ]
16 )
```

| Valid Data augmentation |
|---|

```
1  transform_val = trns.Compose(
2      [
3          trns.Resize((232, 232), interpolation=trns.InterpolationMode.BICUBIC),
4          trns.CenterCrop((224, 224)),
5          trns.ToTensor(),
6          trns.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
7      ]
8  )
```
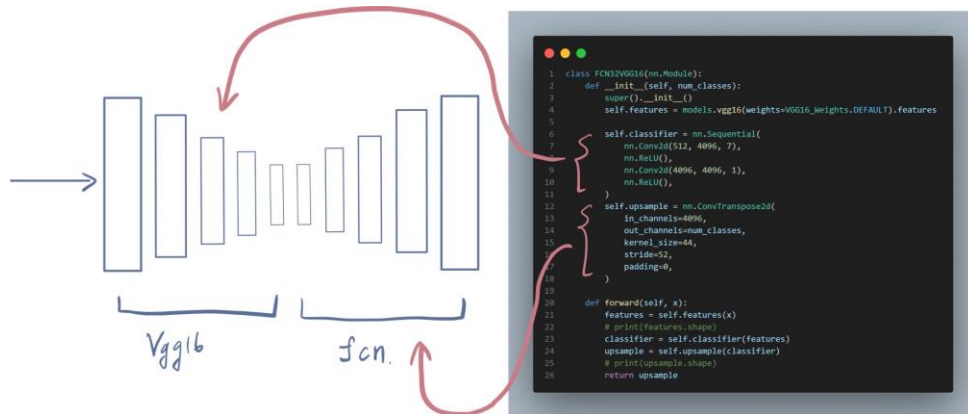
2. *Please conduct the Image classification on Office-Home dataset as the downstream task. Also, please complete the following Table, which contains different image classification setting, and discuss/analyze the results.*

| SSL Info | | | | |
|---|---|---|---|---|
| Learning Rate | Batch Size | Optimizer | Epochs | Backbone |
| 0.0005 | 128 | Adam | 10 | Resnet50 |
| Train Data augmentation | | | | |

```
1  basic_trans = trns.Compose(
2      [
3          trns.Resize([128, 128]),
4          trns.RandomHorizontalFlip(),
5          trns.TrivialAugmentWide(),
6          trns.ToTensor(),
7          trns.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
8      ]
9  )
```

| Fine Tuning Info (ABCDE 都用此規格) | | | | |
|---|---|---|---|---|
| Learning Rate | Batch Size | Optimizer | Epochs | Loss Func |
| 0.0005 | 64 | Adam | 100 | Cross Entropy |
| Train Data augmentation | | | | |

```
basic_transform = [
    # additional data argument
    trns.RandomResizedCrop(224, scale=(0.08, 1.0), ratio=(3.0 / 4.0, 4.0 / 3.0)),
    trns.RandomHorizontalFlip(),  # 水平翻轉
    trns.ColorJitter(brightness=0.5, contrast=0.5, hue=0.5),
]

transform_train = trns.Compose(
    [
        trns.Resize((232, 232), interpolation=trns.InterpolationMode.BICUBIC),
        trns.CenterCrop((224, 224)),
        trns.RandomChoice(basic_transform),
        trns.ToTensor(),
        trns.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
    ]
)
```

| Valid Data augmentation |
| --- |

```
transform_val = trns.Compose(
    [
        trns.Resize((232, 232), interpolation=trns.InterpolationMode.BICUBIC),
        trns.CenterCrop((224, 224)),
        trns.ToTensor(),
        trns.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
    ]
)
```

| Discuss/Analyze the results | | | | |
| --- | --- | --- | --- | --- |
| A | B | C | D | E |
| 41% | 59.5% | 49.6% | 53.9% | 34% |

- A 為單純用 ResNet50 訓練整個 model 的結果，結果比 E 只訓練 Classifier 的好。
- B 是用助教提供的 pre-train weight 去訓練整個 ResNet50，結果優於我寫的 SSL (BYOL)，也是所有選項中最高的。
- 我用 BYOL 弄出的 C 的結果比助教提供的 pre-train weight 的 B 與 D 都低，但都比沒用 pre-train weight 訓練的模型高。
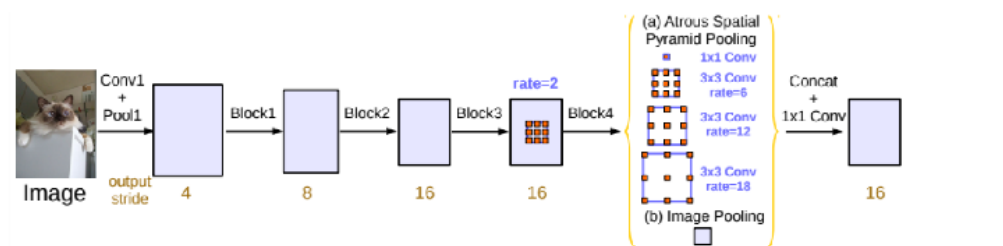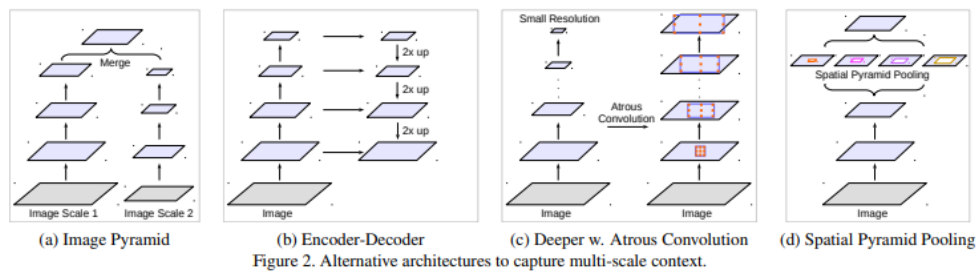- D 與 E 皆比 B 與 C 還低，代表單純只訓練 classifier 是完全不夠的，觀察 E 與 B 的差距可見，我的 SSL(BYOL)提供的效果有限。

● **Problem 3: Semantic Segmentation Task Definition**

*1. Draw the network architecture of your VGG16-FCN32s model (model A).*



*2. Draw the network architecture of the improved model (model B) and explain it differs from your VGG16-FCN32s model.*
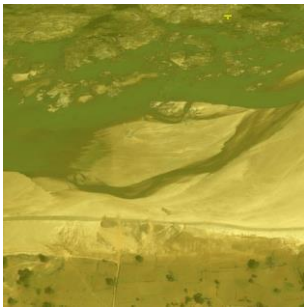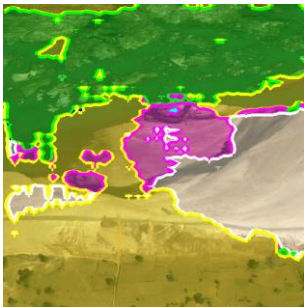
B model 我用 Troch Vision 的 DeepLab3_Resnet50 來 train



(a) Image Pyramid (b) Encoder-Decoder (c) Deeper w. Atrous Convolution (d) Spatial Pyramid Pooling
Figure 2. Alternative architectures to capture multi-scale context.



(原 paper 的架構)

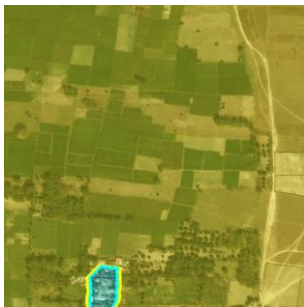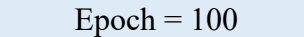B 跟 A 的架構大體上是差不多的,一個 Encoder 與一個 Decoder,差別不同在於他們處理 Layer 的方式不一樣。A 的 Encoder 是用 Vgg16,B 則是用 Resnet50。

DeepLabv3 使用了空洞捲積(Dilated Convolution)作為其主要特點,這有助於捕獲多尺度的語義資訊。它還包括了空間金字塔池化(ASPP)模組,用於在不同尺度上捕獲上下文資訊。DeepLabv3 基於其先進的架構和多尺度上下文資訊捕獲,通常能夠獲得更好的語義分割性能。 它在許多分割任務中表現出色,特別是在複雜場景和小目標分割方面。
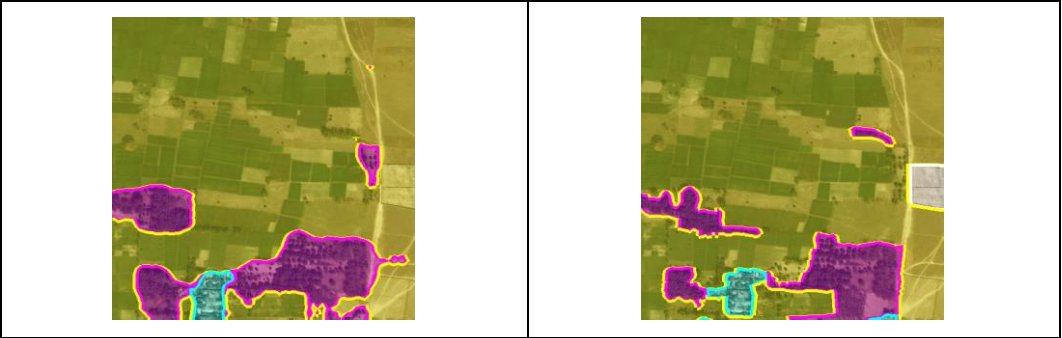
3. *Report mIoUs of two models on the validation set.*

| VGG16-FCN32 mIoU | DeepLabv3_Resnet50 mIoU |
|:---:|:---:|
| 56.4% | 72.9% |

4. *Show the predicted segmentation mask of "validation/0013_sat.jpg" ,*
   *"validation/0062_sat.jpg", "validation/0104_sat.jpg" during the early, middle,*
   *and the final stage during the training process of the improved model.*

| validation/0013_sat.jpg | |
|:---:|:---:|
| Epoch = 1 | Epoch = 50 |
|  |  |
| Epoch = 100 | Ans |
|  |  |

| validation/0062_sat.jpg | |
|:---:|:---:|
| Epoch = 1 | Epoch = 50 |
|  |  |
| Epoch = 100 | Ans |

| validation/0104_sat.jpg | |
| --- | --- |
| Epoch = 1 | Epoch = 50 |
|  |  |
| Epoch = 100 | Ans |
|  |  |