

Deep Learning for Computer Vision

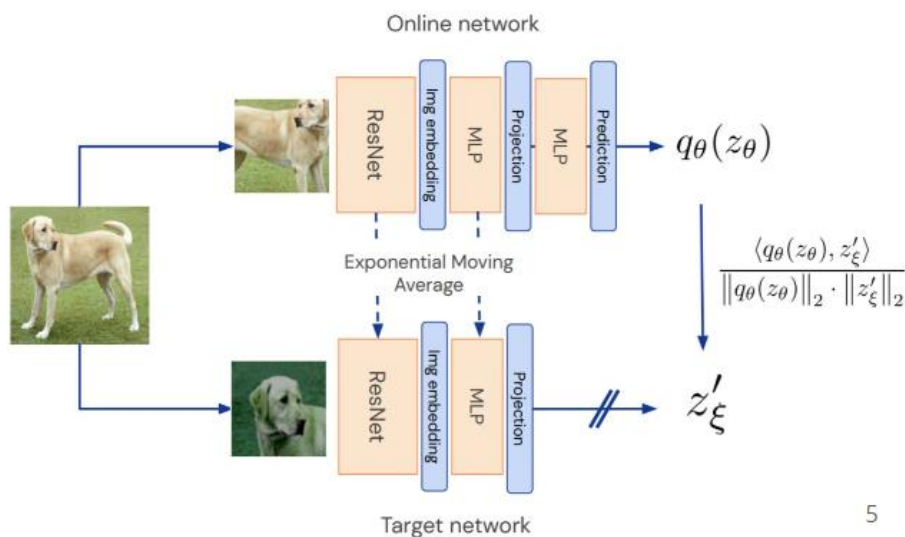
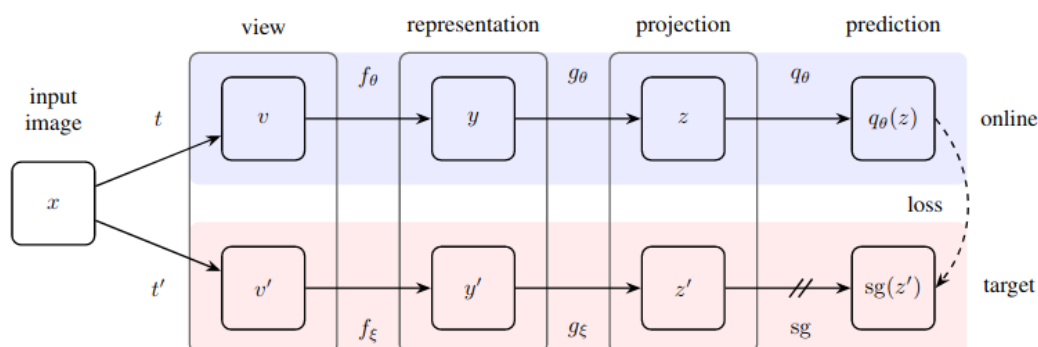
HW1

電子所 ICS 組, R13943015, 張根齊

Problem 1: Self-supervised pre-training for image classification

1.(5%) Describe the implementation details of your SSL method for pre-training the ResNet50 backbone. (including but not limited to the name of the SSL method & data augmentation techniques you used, learning rate schedule, optimizer, and batch size setting for this pre-training phase)

(1) SSL method: BYOL



1. First, apply different data augmentation, t and t' , to the input image x , resulting two views, v and v' .
2. Feed augmented view v and v' into online network Resnet f_θ and target network Resnet f_ε . Output are representation y and y' , respectively.
3. Pass the representation y and y' to online MLP g_θ and target MLP g_ε . Output are projection z and z' respectively.
4. Apply a predictor q_θ to the projection z , and get the prediction $q_\theta(z)$
5. The loss above is computed between $q_\theta(z)$ and z' .

(2) Pre-training setting

1. data augmentation:

`train_mean` and `train_std` are mean and std of ImageNet.

```
# Data Transformations
transform_train = transforms.Compose([
    transforms.Resize(128),
    transforms.CenterCrop(128),
    transforms.RandomHorizontalFlip(),
    transforms.RandomRotation(degrees=10),
    transforms.ToTensor(),
    transforms.Normalize(mean=train_mean, std=train_std),
])
```

2. learning rate schedule and optimizer:

I use Adam as my optimizer. The initial learning rate is 0.001, which is halved every 20 epochs.

```
# Initialize optimizer
optimizer = optim.Adam(learner.parameters(), lr=0.001)
# Different base learning rate and update strategy #!!!!!!
scheduler = torch.optim.lr_scheduler.StepLR(optimizer, step_size=20, gamma=0.5, last_epoch=-1)
```

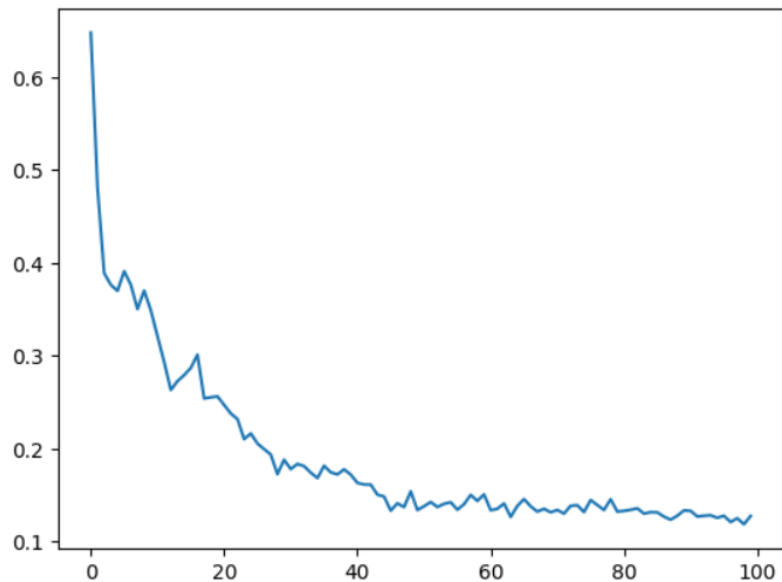
3. batch size and training epoch:

Batch size = 32. Number of epoch = 100.

4. model saving:

The model with the lowest loss is saved.

(3) Pre-training loss over epochs:



2.(20%) 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.

(1) Training setting:

1. data augmentation:

train_mean and train_std are mean and std of ImageNet.

```
# Data Transformations
transform_train = transforms.Compose([
    transforms.Resize(128),
    transforms.CenterCrop(128),
    transforms.RandomHorizontalFlip(),
    transforms.RandomRotation(degrees=10),
    transforms.ColorJitter(brightness=0.2, contrast=0.2, saturation=0.2, hue=0.1),
    transforms.RandomPerspective(distortion_scale=0.5, p=0.5),
    transforms.GaussianBlur(kernel_size=(5, 5), sigma=(0.1, 2.0)),
    transforms.ToTensor(),
    transforms.Normalize(mean=train_mean, std=train_std),
    transforms.RandomErasing(p=0.3, scale=(0.02, 0.05), ratio=(0.3, 3.3), value=0)
])

transform_val = transforms.Compose([
    transforms.Resize(128),
    transforms.CenterCrop(128),
    transforms.ToTensor(),
    transforms.Normalize(mean=train_mean, std=train_std),
])
```

2. learning rate schedule and optimizer:

I use Adam as my optimizer. The initial learning rate is 0.001, which is halved if validation accuracy does not improve for 5 epochs. Also, I applied weight decay to prevent overfitting.

```
# Initialize optimizer
optimizer = optim.Adam(net.parameters(), lr=0.001, weight_decay=1e-4)
# Different base learning rate and update strategy #!!!!!!
scheduler = ReduceLROnPlateau(optimizer, mode='max', factor=0.5, patience=5) #if 5個epoch沒進步==> lr *= 0.5
```

3. Criterion:

I use label smoothing loss to reduce overfitting.

```
# Create Label Smoothing loss instance
criterion = LabelSmoothingLoss(num_classes=65, smoothing=0.1)
```

4. batch size and training epoch:

Batch size = 128. Number of epoch = 100.

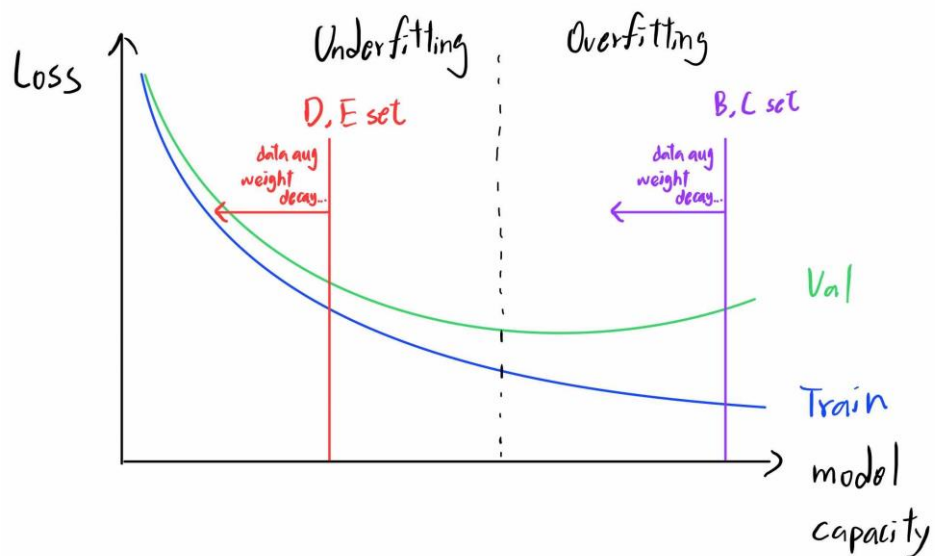
5. model saving:

The model with the highest validation accuracy is saved.

(2) Accuracy of 5 sets:

Setting	Pre-training (Mini-ImageNet)	Fine-tuning (Office-Home)	Validation accuracy (Office-Home)
A	-	Train full model (backbone + classifier)	49.01%
B	w/ label (TAs have provided this backbone)	Train full model (backbone + classifier)	52.71%
C	w/o label (Your SSL pre-trained backbone)	Train full model (backbone + classifier)	46.31%
D	w/ label (TAs have provided this backbone)	Fix the backbone. Train classifier only	37.19%
E	w/o label (Your SSL pre-trained backbone)	Fix the backbone. Train classifier only	8.62%

1. Compare B set with C set, we can see that SL-pretrained backbone outperforms SSL-pretrained backbone. A similar observation can be made when comparing the D set with the E set.
2. Compare {B, C} set with {D, E} set, we can see that {B, C} set outperforms {D, E} set. I guess this is because I applied many techniques to reduce overfitting, ex: 6 data augmentations, label smoothing loss and weight decay. This may cause D and E set do not train well due to their low model capacity(only classifier weights can be updated).

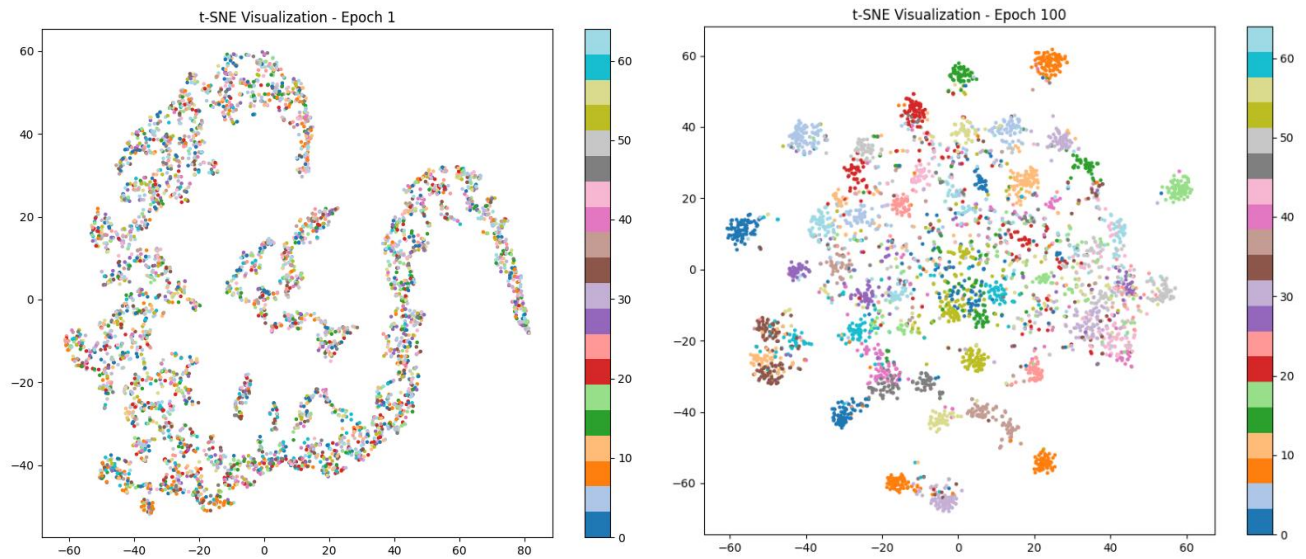


3. We can notice that A set outperforms C set. However, the validation accuracy of A set is 43% at first, since I forget to save the training log so I train it again and result in 49% accuracy. I think this is because SSL backbone does not perform well, and also techniques used to prevent overfitting may cause variations in the training results.

3.(5%) Visualize the learned visual representation of setting C on the train set by implementing t-SNE (t-distributed Stochastic Neighbor Embedding) on the output of the second last layer. Depict your visualization from both the first and the last epochs. Briefly explain the results.

1. TSNE from epoch 1 and epoch 100:

We can see that the points in the t-SNE plot at epoch 1 are mixed and have an irregular distribution. However, in the last epoch, the points of the same color are clustered together and have symmetrical distribution.



Problem 2: Semantic segmentation

1.(3%) Draw the network architecture of your VGG16-FCN32s model (model A).

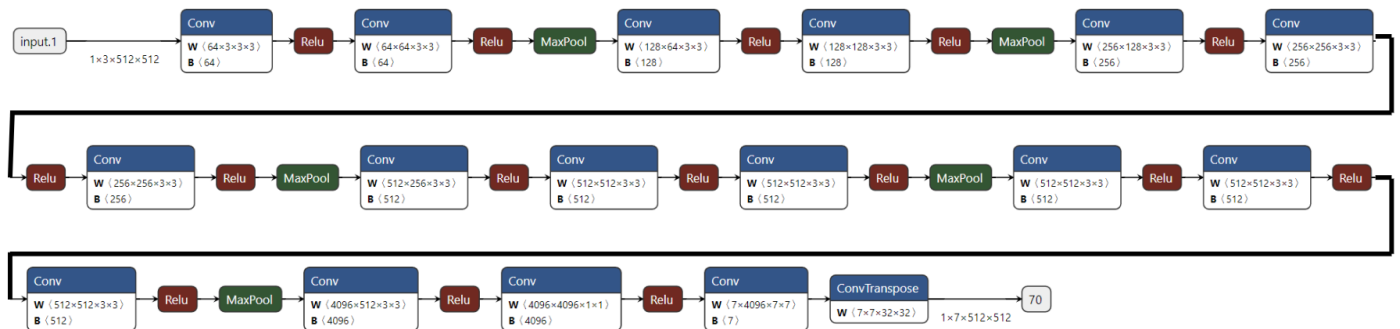
I replace FC layers of VGG16 by Convolutional layers, which has same size as output feature map of Maxpool2d-31 ($H = 16, W = 16$) and has same #channel as original FC layers ($4096 \rightarrow 4096 \rightarrow 7$). Since FCN32s need to up-sample the feature map by 32x back to the input image size, so I add a transpose convolution layer to up-sample $7 * 16 * 16$ feature map to $7 * 512 * 512$.

Layer	(type)	Output Shape	Param #
1			
2			
3			
4	Conv2d-1	$[-1, 64, 512, 512]$	1,792
5	ReLU-2	$[-1, 64, 512, 512]$	0
6	Conv2d-3	$[-1, 64, 512, 512]$	36,928
7	ReLU-4	$[-1, 64, 512, 512]$	0
8	MaxPool2d-5	$[-1, 64, 256, 256]$	0
9	Conv2d-6	$[-1, 128, 256, 256]$	73,856
10	ReLU-7	$[-1, 128, 256, 256]$	0
11	Conv2d-8	$[-1, 128, 256, 256]$	147,584
12	ReLU-9	$[-1, 128, 256, 256]$	0
13	MaxPool2d-10	$[-1, 128, 128, 128]$	0
14	Conv2d-11	$[-1, 256, 128, 128]$	295,168
15	ReLU-12	$[-1, 256, 128, 128]$	0
16	Conv2d-13	$[-1, 256, 128, 128]$	590,080
17	ReLU-14	$[-1, 256, 128, 128]$	0
18	Conv2d-15	$[-1, 256, 128, 128]$	590,080
19	ReLU-16	$[-1, 256, 128, 128]$	0
20	MaxPool2d-17	$[-1, 256, 64, 64]$	0
21	Conv2d-18	$[-1, 512, 64, 64]$	1,180,160
22	ReLU-19	$[-1, 512, 64, 64]$	0
23	Conv2d-20	$[-1, 512, 64, 64]$	2,359,808
24	ReLU-21	$[-1, 512, 64, 64]$	0
25	Conv2d-22	$[-1, 512, 64, 64]$	2,359,808
26	ReLU-23	$[-1, 512, 64, 64]$	0
27	MaxPool2d-24	$[-1, 512, 32, 32]$	0
28	Conv2d-25	$[-1, 512, 32, 32]$	2,359,808
29	ReLU-26	$[-1, 512, 32, 32]$	0
30	Conv2d-27	$[-1, 512, 32, 32]$	2,359,808
31	ReLU-28	$[-1, 512, 32, 32]$	0
32	Conv2d-29	$[-1, 512, 32, 32]$	2,359,808
33	ReLU-30	$[-1, 512, 32, 32]$	0
34	MaxPool2d-31	$[-1, 512, 16, 16]$	0
35	Conv2d-32	$[-1, 4096, 16, 16]$	18,878,464
36	ReLU-33	$[-1, 4096, 16, 16]$	0
37	Dropout2d-34	$[-1, 4096, 16, 16]$	0
38	Conv2d-35	$[-1, 4096, 16, 16]$	16,781,312

```

39      ReLU-36      [-1, 4096, 16, 16]      0
40      Dropout2d-37      [-1, 4096, 16, 16]      0
41      Conv2d-38      [-1, 7, 16, 16]      1,404,935
42      ConvTranspose2d-39      [-1, 7, 512, 512]      50,176
43      -----
44      Total params: 51,829,575
45      Trainable params: 51,829,575
46      Non-trainable params: 0
47      -----
48      Input size (MB): 3.00
49      Forward/backward pass size (MB): 1203.01
50      Params size (MB): 197.71
51      Estimated Total Size (MB): 1403.73
52      -----
53

```

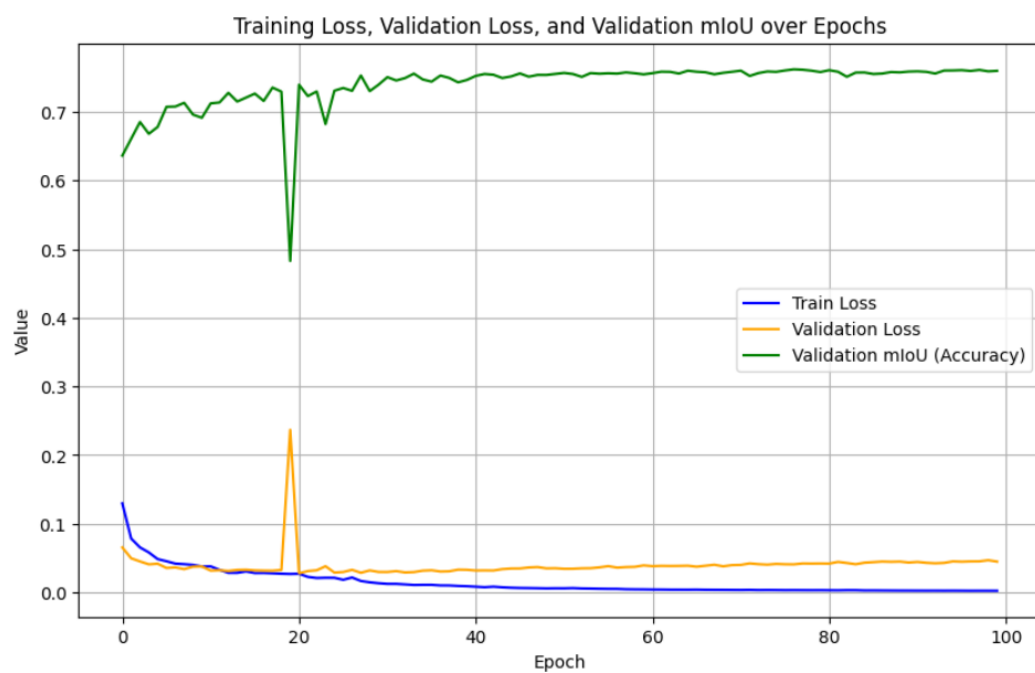


2.(3%) Draw the network architecture of the improved model (model B)

and explain it differs from your VGG16-FCN32s model.

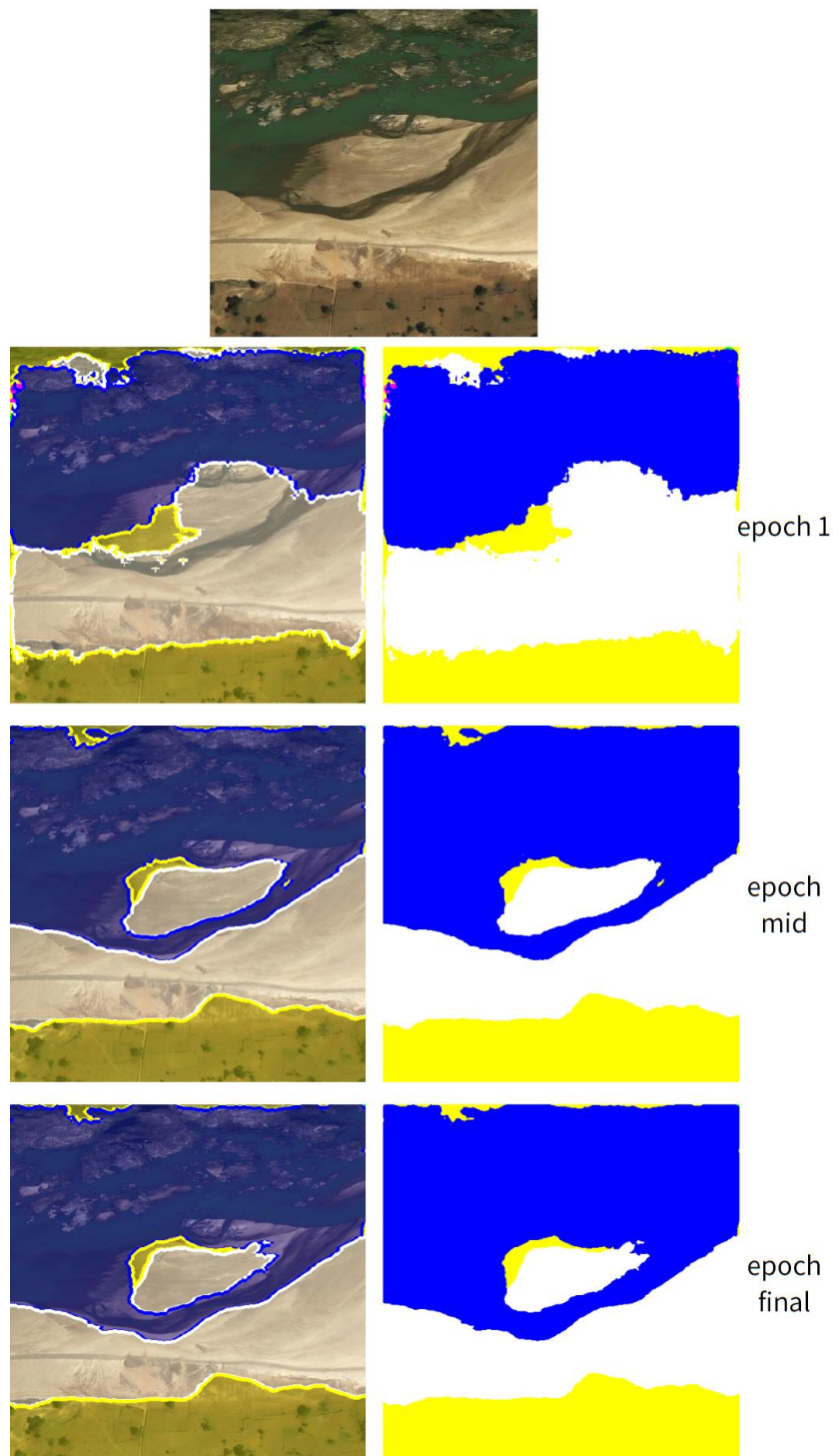
I use DeepLabV3 on resnet101 as my B model. First, they use different backbone, resnet101 v.s. VGG16. Second, DeepLabV3 introduced atrous(dilated) convolution and the ASPP module applies multiple parallel atrous convolutions with different dilation rates, which allows model to capture multi-scale features without losing resolution.

Training loss, validation loss, mIoU over epochs:



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.

(1) 0013_sat:



(2) 0062_sat:



epoch 1

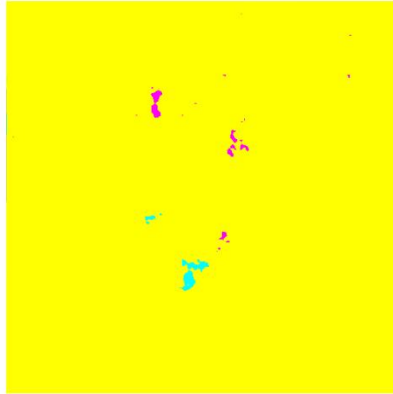


epoch
mid

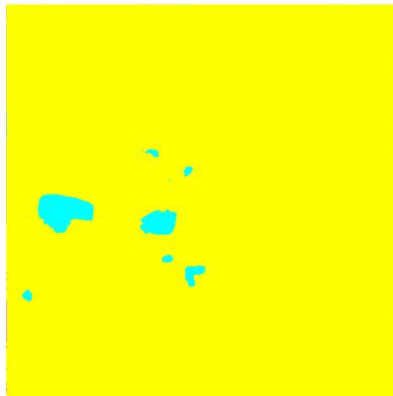


epoch
final

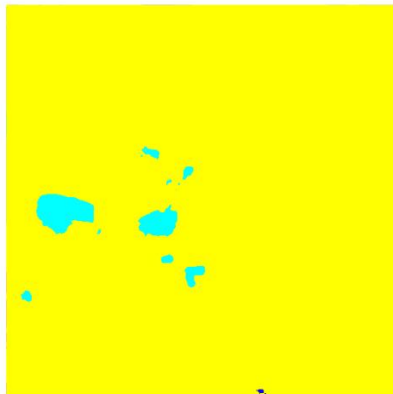
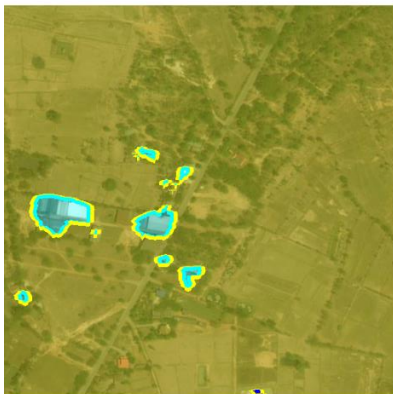
(3) 0104_sat:



epoch 1



epoch
mid



epoch
final

5.(10%) Use segment anything model (SAM) to segment three of the images in the validation dataset, report the result images and the method you use.

I use meta demo website(<https://segment-anything.com/demo>) to inference SAM on 0055_sat, 0061_sat, 0101_sat. I use ViT-H as the pre-trained image encoder and prompt SAM with each point in the grid.



Reference:

1. Chatgpt:

<https://openai.com/chatgpt/>

2. Bootstrap Your Own Latent (BYOL), in Pytorch:

<https://github.com/lucidrains/byol-pytorch>

3. CosineAnnealingLR — PyTorch 2.4 documentation:

https://pytorch.org/docs/stable/generated/torch.optim.lr_scheduler.CosineAnnealingLR.html

4. DeepLabV3 — Torchvision 0.19 documentation

<https://pytorch.org/vision/stable/models/deeplabv3.html>

5. Semantic Segmentation on PASCAL VOC 2012 test:

<https://paperswithcode.com/sota/semantic-segmentation-on-pascal-voc-2012>

6. Segment Anything-Research by Meta AI:

<https://segment-anything.com/demo>