

You can modify this report template, and upload your results in **PDF format**. Reports in other forms/formats will result in **ZERO point**. Reports written in either Chinese or English are both acceptable. The length of your report should **NOT** exceed **6 pages (excluding bonus)**.

Name: 謝子然 Dep.:資工碩一 Student ID:R07922089

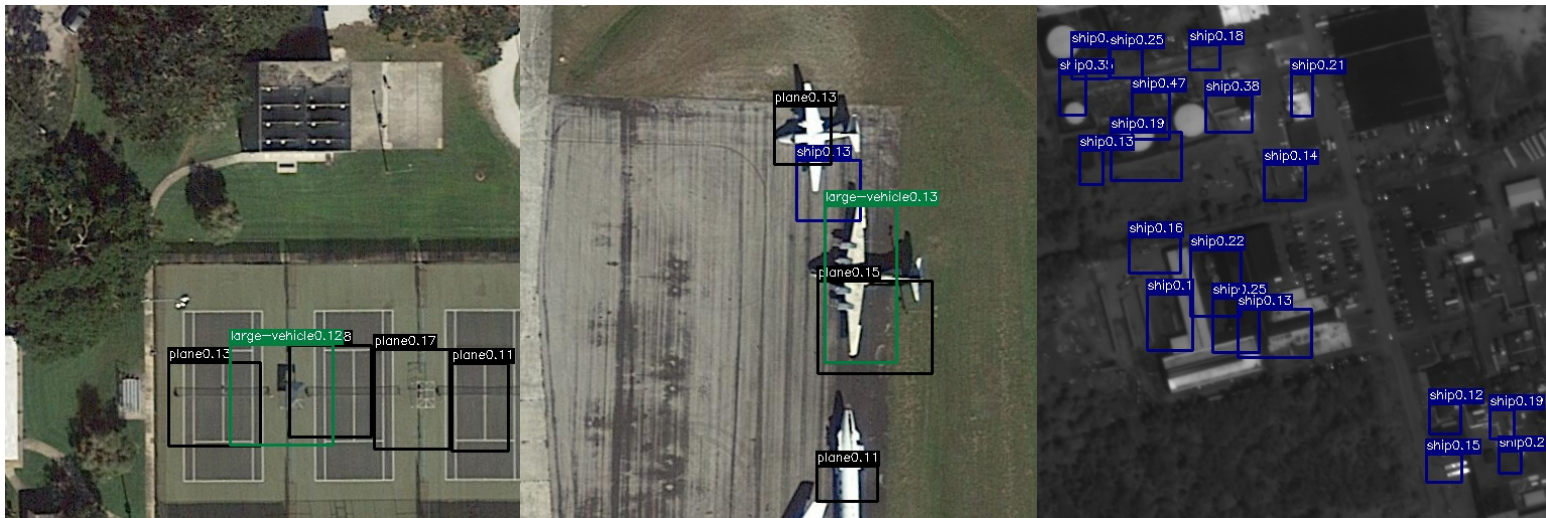
1. (5%) Print the network architecture of your YoloV1-vgg16bn model and describe your training config. (optimizer,batch size....and so on)

```
13 class VGG(nn.Module):
14
15     def __init__(self, features, output_size=1274, image_size=448):
16         # output: 1274 is 7 x 7 x (2 x 5 + 16)
17         super(VGG, self).__init__()
18         self.features = features
19         self.image_size = image_size
20
21         self.yolo = nn.Sequential(
22             # TODO
23             nn.Linear(512 * 7 * 7, 4096),
24             nn.ReLU(True),
25             nn.Dropout(),
26             nn.Linear(4096, output_size)
27         )
28         self._initialize_weights()
29
30     def forward(self, x):
31         # [batch, 3, 488, 488]
32         x = self.features(x)
33         # [1, 512, 7, 7]
34         x = x.view(x.size(0), -1)
35         # [1, 25088]
36         x = self.yolo(x)
37         x = torch.sigmoid(x)
38         x = x.view(-1, 7, 7, 26) # note the position of "depth(26)" is different
39         return x
40
41     def _initialize_weights(self):
42         for m in self.modules():
43             if isinstance(m, nn.Linear):
44                 m.weight.data.normal_(0, 0.01)
45                 m.bias.data.zero_()
46
47
48 def make_layers(cfg, batch_norm=False):
49     layers = []
50     in_channels = 3
51     s = 1
52     first_flag=True
53     for v in cfg:
54         s=1
55         if (v==64 and first_flag):
56             # for reducing feature map size
57             s=2
58             first_flag=False
59         if v == 'M':
60             layers += [nn.MaxPool2d(kernel_size=2, stride=2)]
61         else:
62             conv2d = nn.Conv2d(in_channels, v, kernel_size=3, stride=s, padding=1)
63             if batch_norm:
64                 layers += [conv2d, nn.BatchNorm2d(v), nn.ReLU(inplace=True)]
65             else:
66                 layers += [conv2d, nn.ReLU(inplace=True)]
67             in_channels = v
68     return nn.Sequential(*layers)
69
70 def conv_bn_relu(in_channels,out_channels,kernel_size=3,stride=2,padding=1):
71     return nn.Sequential(
72         nn.Conv2d(in_channels,out_channels,kernel_size=kernel_size,padding=padding,stride=stride),
73         nn.BatchNorm2d(out_channels),
74         nn.ReLU(True)
75     )
76
77
78 cfg = {
79     'D': [64, 64, 'M', 128, 128, 'M', 256, 256, 256, 'M', 512, 512, 512, 'M', 512, 512, 512, 'M'],
80 }
81
82
83
84
85 def YoloV1_vgg16bn(pretrained=False, **kwargs):
86     """VGG 16-layer model (configuration "D") with batch normalization
87     Args:
88         pretrained (bool): If True, returns a model pre-trained on ImageNet
89     """
90     yolo = VGG(make_layers(cfg['D'], batch_norm=True), **kwargs)
91     if pretrained:
92         vgg_state_dict = model_zoo.load_url(model_urls['vgg16_bn'])
93         yolo_state_dict = yolo.state_dict()
94         for k in vgg_state_dict.keys():
95             if k in yolo_state_dict.keys() and k.startswith('features'):
96                 yolo_state_dict[k] = vgg_state_dict[k]
97         yolo.load_state_dict(yolo_state_dict)
98     return yolo
99
100
101 models.py 32,28 21% models.py 75,5 78%
```

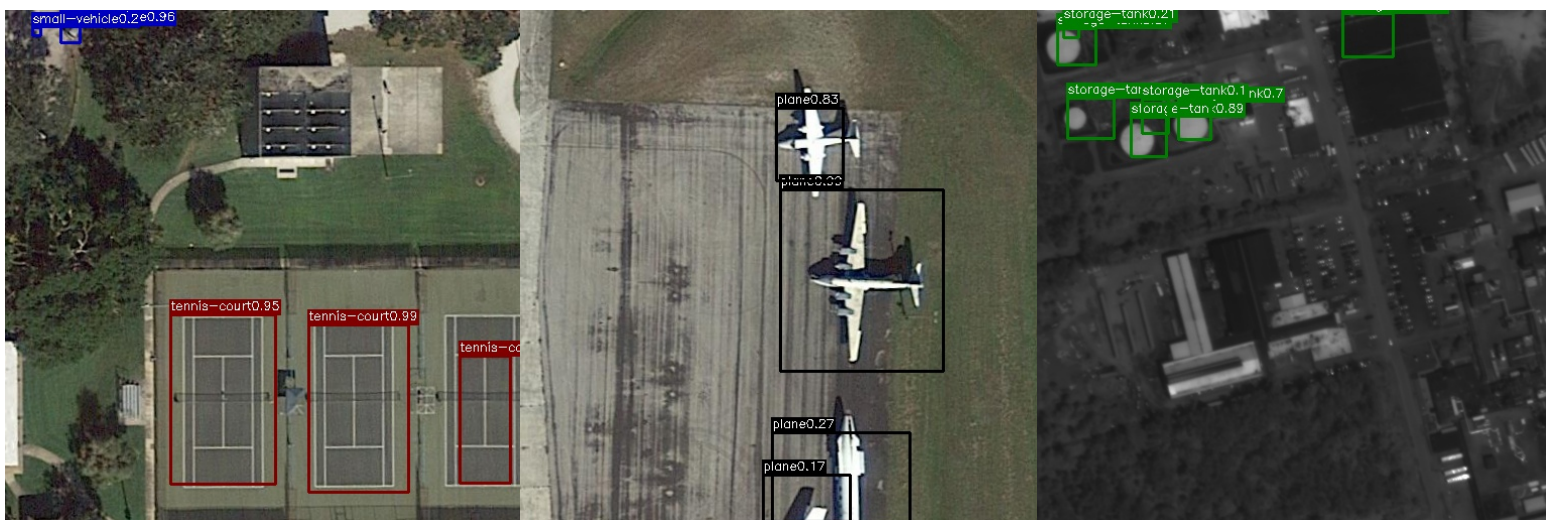
batch size = 15, learning rate fix at 10^{-3} , lambda_coord = 5, lambda_noobj=0.5, optimizer = SGD with momentum = 0.9, weight decay = $5e^{-4}$, max epoch = 100

2. (10%) Show the predicted bbox image of “val1500/0076.jpg”, “val1500/0086.jpg”, “val1500/0907.jpg” during the early, middle, and the final stage during the training stage. (For example, results of 1st, 10th, 20th epoch)

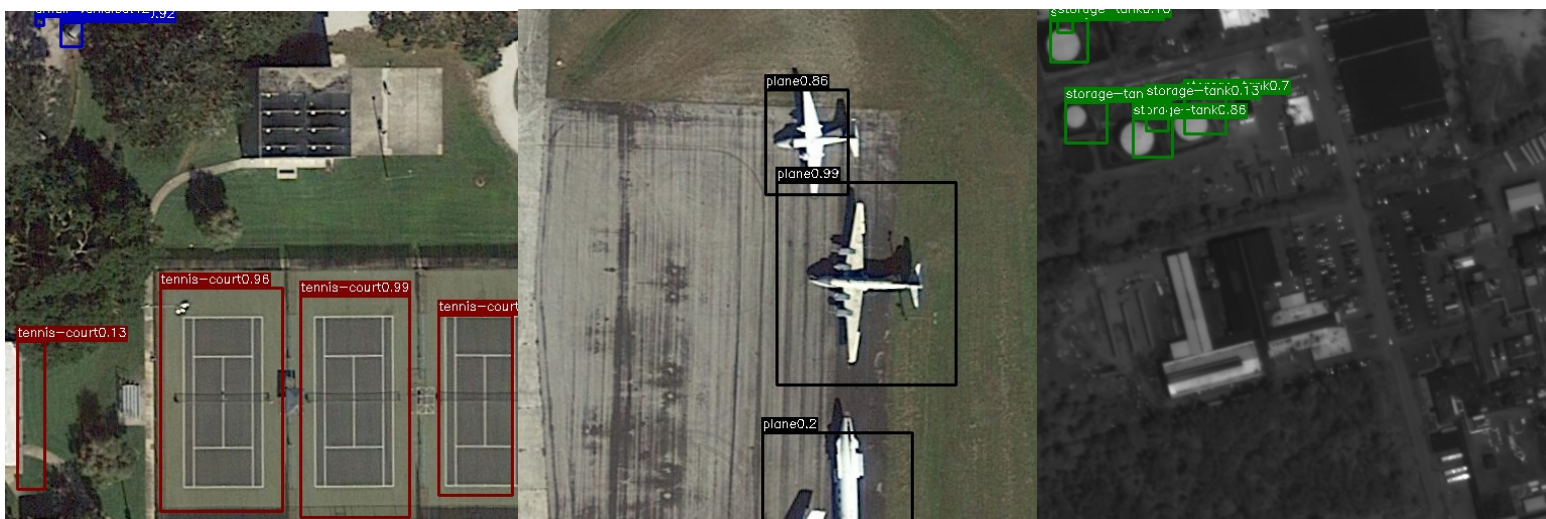
epoch 1



epoch 38



epoch 58



3. (10%) Implement an improved model which performs better than your baseline model. Print the network architecture of this model and describe it.

```

1 import torch.nn as nn
2 import math
3 import torch
4 import torch.utils.model_zoo as model_zoo
5
6
7 __all__ = ['vgg16_bn', 'vgg19_bn']
8 model_urls = {
9     'vgg16_bn': 'https://download.pytorch.org/models/vgg16_bn-6c64b313.pth',
10    'vgg19_bn': 'https://download.pytorch.org/models/vgg19_bn-c79401a0.pth',
11 }
12
13
14 class VGG(nn.Module):
15
16     def __init__(self, features, output_size=1274, image_size=448):
17         # output: 1274 is 7 x 7 x (2 x 5 + 16)
18         super(VGG, self).__init__()
19         self.features = features
20         self.image_size = image_size
21
22         self.yolo = nn.Sequential(
23             #TODD
24             nn.Linear(512 * 7 * 7 * 2, 4096),
25             nn.ReLU(True),
26             nn.Dropout(),
27             nn.Linear(4096, output_size)
28         )
29         self._initialize_weights()
30
31     def forward(self, x):
32         # [batch, 3, 488, 488]
33         x = self.features(x)
34         # [1, 512, 7, 7]
35         x = x.view(x.size(0), -1)
36         # [1, 25088]
37         x = self.yolo(x)
38         x = torch.sigmoid(x)
39         x = x.view(-1, 7, 7, 26) # note the position of "depth(26)" is different
40         return x
41
42     def _initialize_weights(self):
43         for m in self.modules():
44             if isinstance(m, nn.Linear):
45                 m.weight.data.normal_(0, 0.01)
46                 m.bias.data.zero_()
47
48
49 def make_layers(cfg, batch_norm=False):
50     layers = []
51     in_channels = 3
52     s = 1
53     first_flag=True
54     for v in cfg:
55         s=1
56         if (v==64 and first_flag):
57             # for reducing feature map size
58             s=2
59             first_flag=False
60         if v == 'M':
61             layers += [nn.MaxPool2d(kernel_size=2, stride=2)]
62         else:
63             conv2d = nn.Conv2d(in_channels, v, kernel_size=3, stride=s, padding=1)
64             if batch_norm:
65                 layers += [conv2d, nn.BatchNorm2d(v), nn.ReLU(inplace=True)]
66             else:
67                 layers += [conv2d, nn.ReLU(inplace=True)]
68             in_channels = v
69     return nn.Sequential(*layers)
70
71 def conv_bn_relu(in_channels, out_channels, kernel_size=3, stride=2, padding=1):
72     return nn.Sequential(
73         nn.Conv2d(in_channels, out_channels, kernel_size=kernel_size, padding=padding, stride=stride),
74         nn.BatchNorm2d(out_channels),
75         nn.ReLU(True)
76     )
77
78
79 cfg = {
80     'D': [64, 64, 'M', 128, 128, 'M', 256, 256, 256, 'M', 512, 512, 512, 'M', 512, 512, 512,
81           'M'],
82     'E': [64, 64, 'M', 128, 128, 'M', 256, 256, 256, 256, 'M', 512, 512, 512, 512, 'M',
83           512, 512, 512, 512, 'M'],
84 }
85
86
87 def YoloV1_vgg19bn(pretrained=False, **kwargs):
88     """VGG 19-layer model (configuration "E") with batch normalization
89     Args:
90         pretrained (bool): If True, returns a model pre-trained on ImageNet
91     """
92     yolo = VGG(make_layers(cfg['E'], batch_norm=True), **kwargs)
93     if pretrained:
94         vgg_state_dict = model_zoo.load_url(model_urls['vgg19_bn'])
95         yolo_state_dict = yolo.state_dict()
96         for k in vgg_state_dict.keys():
97             if k in yolo_state_dict.keys() and k.startswith('features'):
98                 yolo_state_dict[k] = vgg_state_dict[k]
99     yolo.load_state_dict(yolo_state_dict)
100     return yolo
101

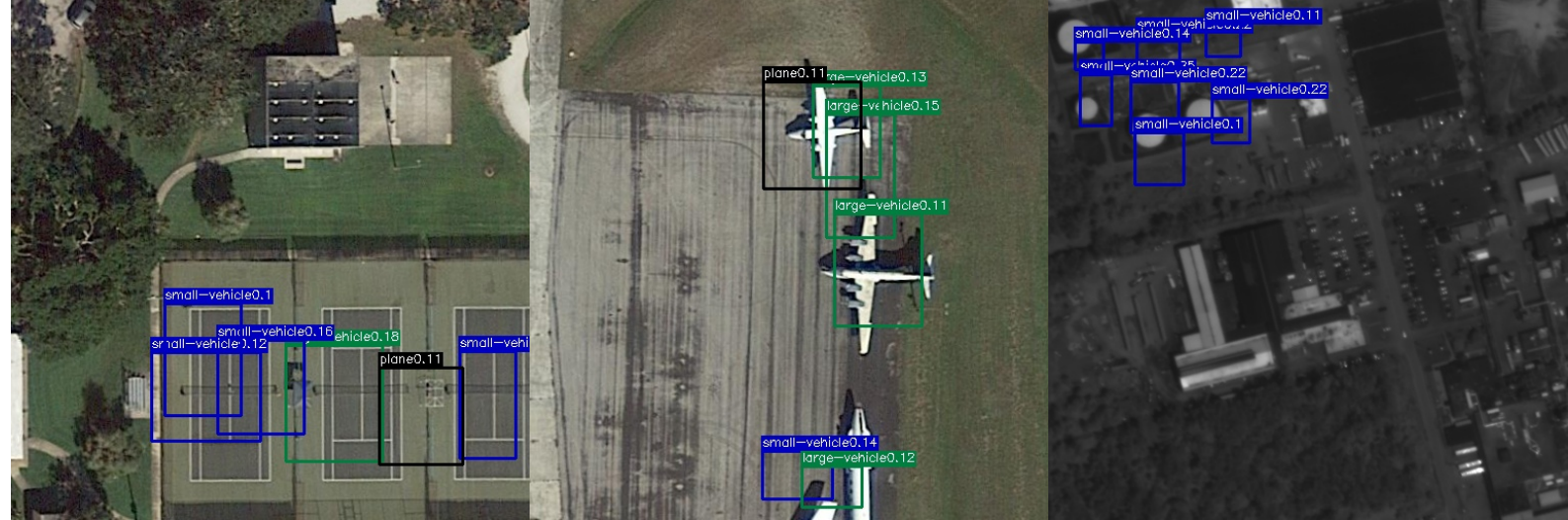
```

I changed the backbone network from VGG16 to VGG19. Besides, and I also change make the learning rate scale down by multiplying 0.8 every 30 epochs.

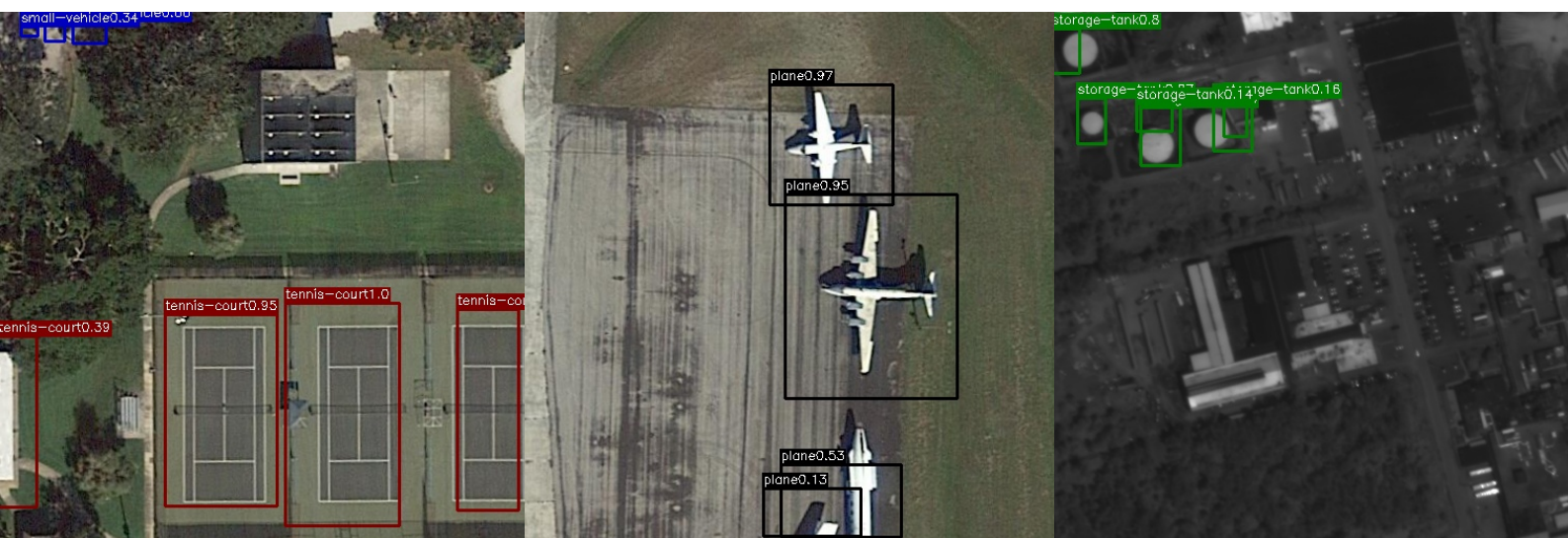
batch size = 15, lambda_coord = 5, lambda_noobj=0.5, optimizer = SGD with momentum = 0.9, weight decay = $5e^{-4}$, max epoch = 100

4. (10%) Show the predicted bbox image of “val1500/0076.jpg”, “val1500/0086.jpg”, “val1500/0907.jpg” during the early, middle, and the final stage during the training process of this improved model.

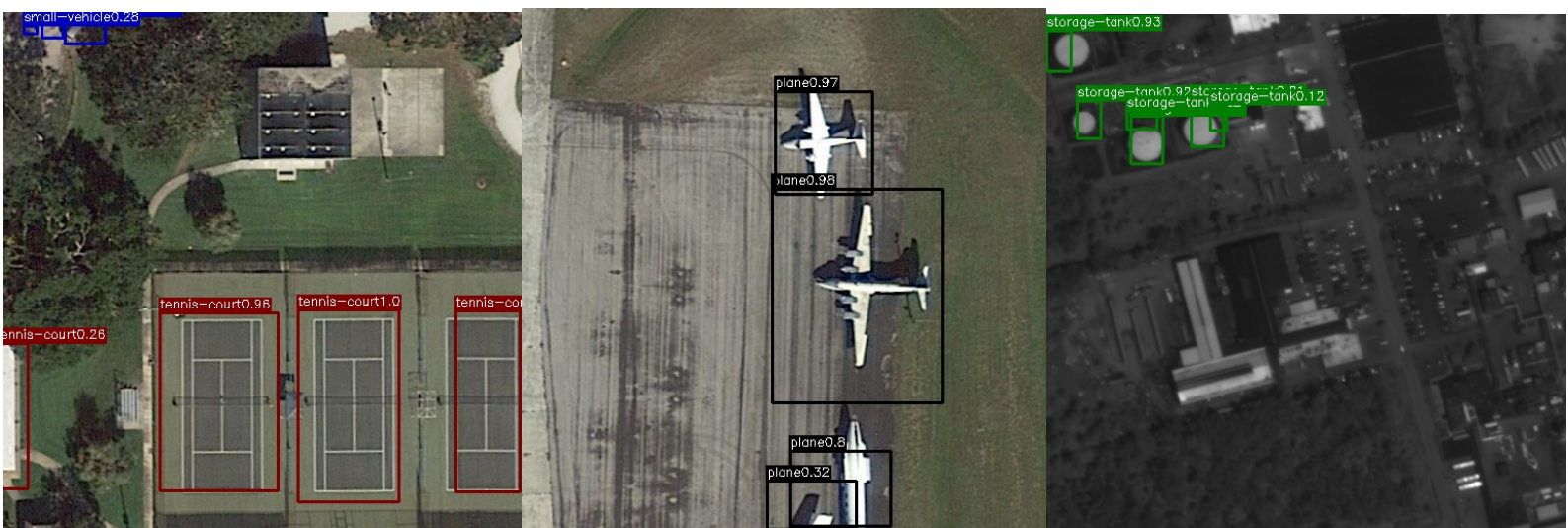
epoch 1



epoch 40



epoch 72



5. (15%) Report mAP score of both models on the validation set. Discuss the reason

why the improved model performs better than the baseline one. You may conduct some experiments and show some evidences to support your reasoning.

mAP of baseline model: 0.13062

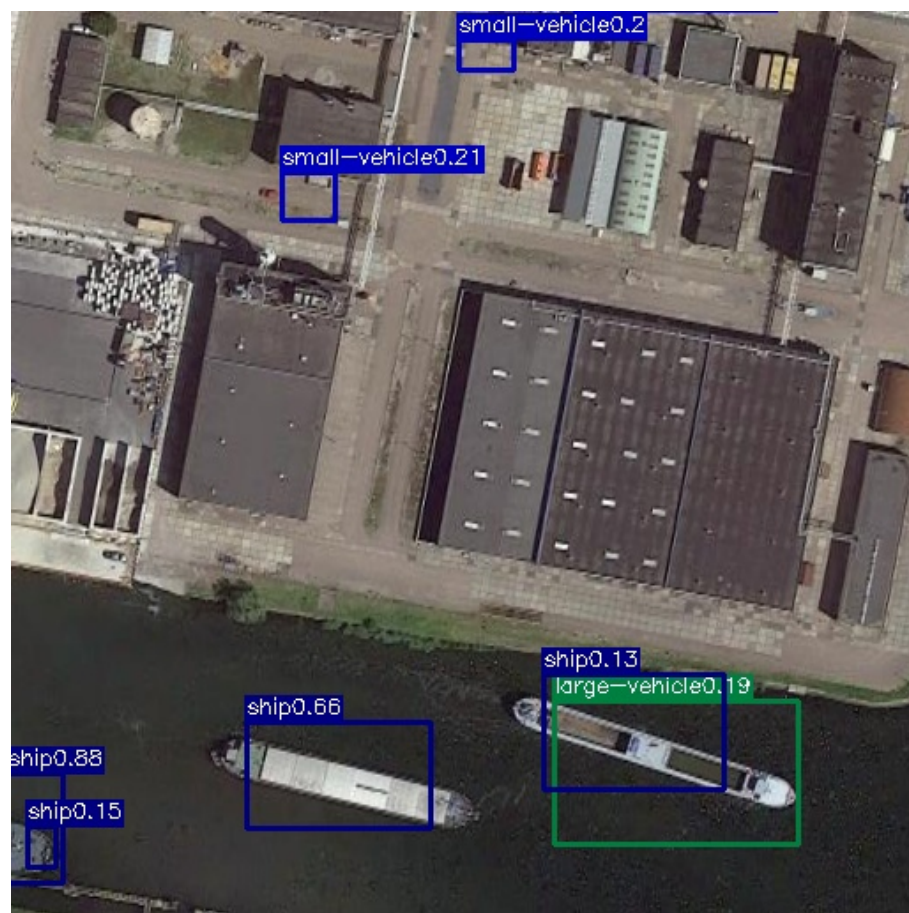
mAP of best model: 0.13217

The baseline model converged at epoch 58, and the best model converged at epoch 72, I think unfixed rate may help the network to find out the true local minimum. Besides, comparing to VGG16, VGG19 may have more power to extract features from input image, as we can see in the comparison of 0086.jpg, my best model do better when the plane is incomplete in image. I conclude that it learned more plane features therefore it has more confidence when it sees incomplete airplanes.

6. **bonus (5%)** Which classes prediction perform worse than others? Why? You should describe and analyze it.

classname: plane	classname: small-vehicle	classname: basketball-court	classname: harbor
ap: 0.36414999560461647	ap: 0.005505289865479439	ap: 0.14545454545454545	ap: 0.10730313165095774
classname: baseball-diamond	classname: large-vehicle	classname: storage-tank	classname: swimming-pool
ap: 0.15844155844155847	ap: 0.12410147991543341	ap: 0.03827751196172249	ap: 0.07692307692307693
classname: bridge	classname: ship	classname: soccer-ball-field	classname: helicopter
ap: 0.15729099505695252	ap: 0.03745595344330556	ap: 0.11363636363636365	ap: 0.0
classname: ground-track-field	classname: tennis-court	classname: roundabout	classname: container-crane
ap: 0.11188811188811189	ap: 0.5651810671753822	ap: 0.1090909090909091	ap: 0.0

Ship sometime may be predicted as large-vehicle, as they actually looks pretty likely in aerial photos. Besides, Storage tank may easily infected by shadow. The following image shows a mis-predicted ship and a low-confidence storage tank.



7. Reference: Loss module and NMS are referenced and modified from: https://github.com/bobo0810/AnnotatedNetworkModelGit/tree/master/Yolov1_pytorch