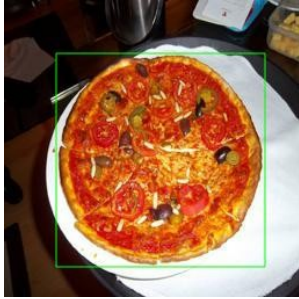


HW5 Report

Cheng-Yun Yang

◆ Creating Your Own Object Localization Dataset

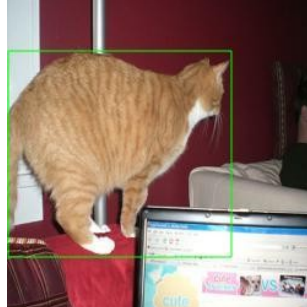
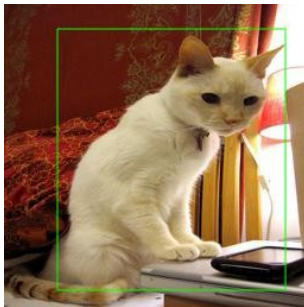
■ Pizza



■ Bus



■ Cat



◆ Building your deep neural network

ResBlock

```
88 class ResBlock(nn.Module):
89     def __init__(self, input_nc, output_nc, stride=1) -> None:
90         super(ResBlock, self).__init__()
91         self.conv1 = nn.Sequential(nn.Conv2d(input_nc, output_nc, kernel_size=3, stride=stride, padding=1),
92                                   nn.BatchNorm2d(output_nc),
93                                   nn.ReLU())
94         self.conv2 = nn.Sequential(nn.Conv2d(output_nc, output_nc, kernel_size=3, stride=1, padding=1),
95                                   nn.BatchNorm2d(output_nc))
96         self.relu = nn.ReLU()
97         # Add a 1x1 conv to match the feature map sizes
98         self.shortcut = nn.Sequential(nn.Conv2d(input_nc, output_nc, kernel_size=1, stride=stride),
99                                       nn.BatchNorm2d(output_nc))
100     def forward(self, x):
101         residual = self.shortcut(x)
102         out = self.conv1(x)
103         out = self.conv2(out)
104         out = out.clone() + residual
105         out = self.relu(out)
106         return out
```

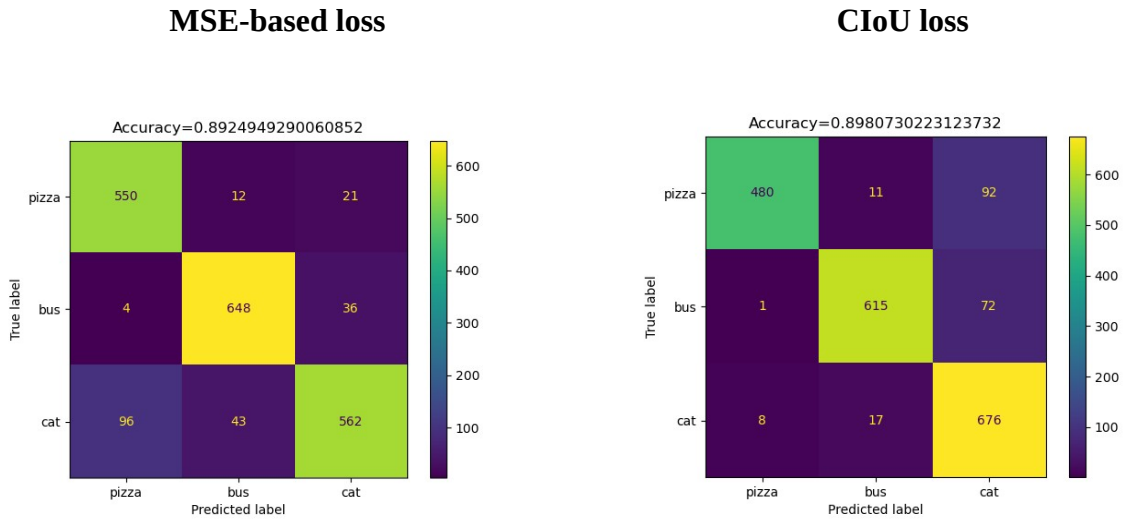
HW5Net

```
108 class HW5Net(nn.Module):
109     def __init__(self, input_nc, ngf=8, n_blocks=4) -> None:
110         assert(n_blocks>=0)
111         super(HW5Net, self).__init__()
112         # The first conv layer
113         model = [nn.ReflectionPad2d(3),
114                 nn.Conv2d(input_nc, ngf, kernel_size=7, padding=0),
115                 nn.BatchNorm2d(ngf),
116                 nn.ReLU(True)]
117         # Add downsampling layers
118         n_downsampling = 4
119         for i in range(n_downsampling):
120             mult = 2**i # 1, 2, 4, 8
121             model += [nn.Conv2d(ngf*mult, ngf*mult*2, kernel_size=3, stride=2, padding=1),
122                     nn.BatchNorm2d(ngf*mult*2),
123                     nn.ReLU(True)]
124         # Add ResNet blocks
125         mult = 2 ** n_downsampling # 16
126         for i in range(n_blocks):
127             model += [ResBlock(ngf*mult, ngf*mult, stride=2)]
128         self.model = nn.Sequential(*model)
129         # Classification head
130         class_head = [nn.Dropout(0.5),
131                      nn.Linear(512,3)]
132         self.class_head = nn.Sequential(*class_head)
133         # Regression head
134         bbox_head = [nn.Dropout(0.5),
135                    nn.Linear(512,64),
136                    nn.ReLU(),
137                    nn.BatchNorm1d(64),
138                    nn.Dropout(0.5),
139                    nn.Linear(64,4)]
140         self.bbox_head = nn.Sequential(*bbox_head)
141     def forward(self, x):
142         ft = self.model(x)
143         ft = ft.view(ft.shape[0], -1) # Change dimension of feature maps to fit linear layers
144         cls = self.class_head(ft)
145         bbox = self.bbox_head(ft)
146         bbox = nn.Sigmoid()(bbox).clone()
147         return cls, bbox
```

- 64 learnable layers in total

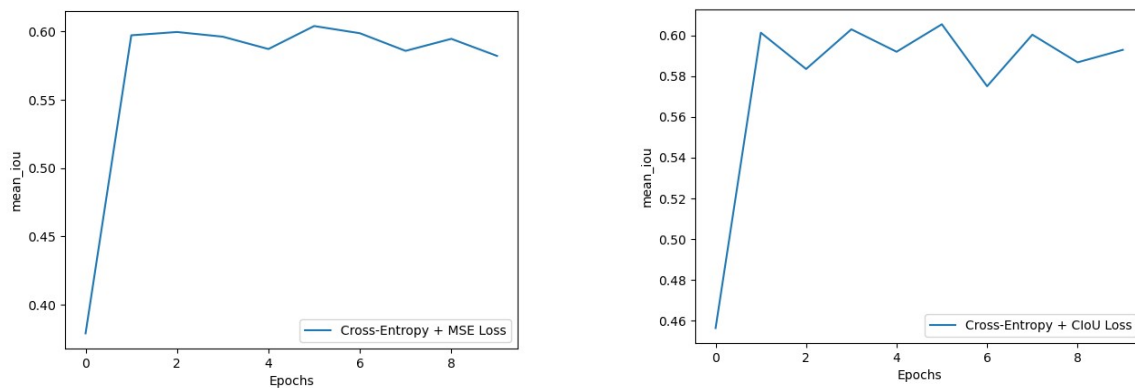
◆ Training and Evaluating Your Trained Network

■ Confusion matrix



Since they both use cross-entropy loss as classification loss, the accuracy is roughly the same.

■ Mean IoU



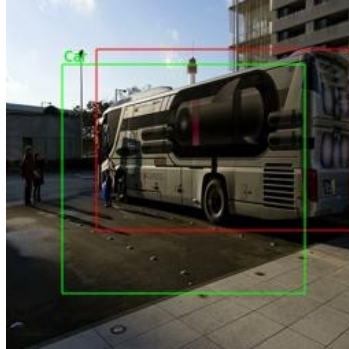
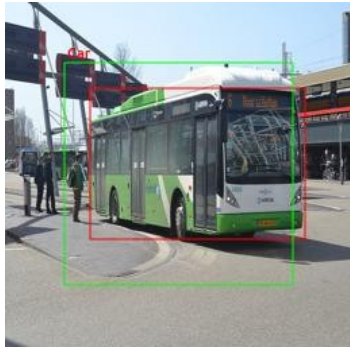
The highest mean IoU over epochs of using CIoU loss is a little bit better than that of MSE loss. Another difference is that training with CIoU is more unstable since it brings larger loss value compared to that of MSE loss.

■ Visualization

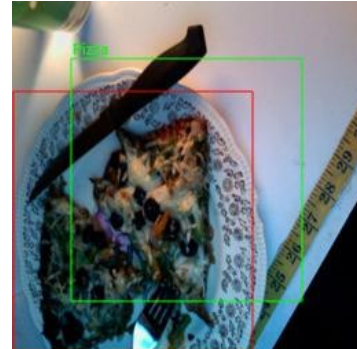
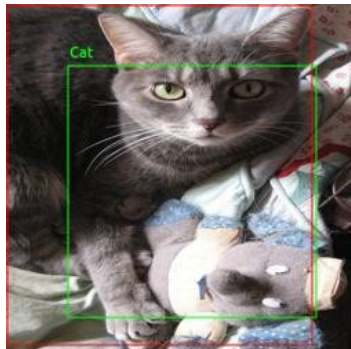
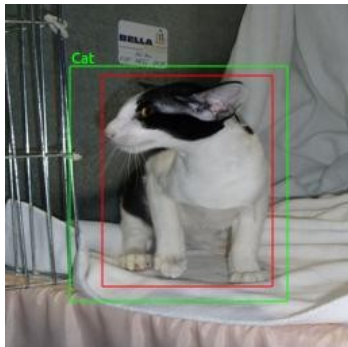
Pizza



Bus



Cat



We can observe that the prediction box tends to locate at the center of the images. The possible reason is that most training samples locate at the the center so the detector also learns this feature. For the right one of the bus sample, we can see the IoU is relatively low because the bus is in the peripheral region of image. This also gives us an intuition of anchor box setting, which we might work in the next assignment. If we assign some pre-defined anchors and learn the offsets to those anchors, it can be expected that the performance will be better. At least we can somewhat limit every bounding box to learn to be a middle normal box. Besides, I observe that the training tends to overfit the training dataset, so if we can do something like early-stopping or add some dropout layers or tune the regularization terms, the problem of overfitting might be alleviated to make the overall performance better.