

# 1 Introduction

The objective of this homework assignment was to be able to apply the YOLO network and skip connections for the purpose of multi-object detection and localization using Microsoft's COCO Dataset.

## 2 Methodology

### 2.1 COCO Downloader

The methodology used to download the images and its respective labels and bounding boxes from the COCO dataset is as follows

1. Get all the image ids that fit all the permutations of the categories (bus, cat, pizza)
2. Iterate through the list of image ids and for each image id create a dictionary containing information on the image id, path to image, and the bounding boxes
3. If the number of annotations are associated to the three categories we care about and the area of the annotation bounding box is greater than 4096, then we add the necessary values to the dictionary created in the previous step
4. The dictionary for each image is then appended to a pandas DataFrame for easy access

The number of training images obtained from the COCO dataset is 6883 and the number of validation images obtained are 3491

#### 2.1.1 Implementation

```
1 # %%
2 from pycocotools.coco import COCO
3 from PIL import Image
4 import numpy as np
5 import os
6 import pandas as pd
7 import torchvision.transforms as tvf
8 from torchvision.io import read_image
9 import argparse
10 from tqdm import tqdm
11 import itertools
12 from pprint import pprint
13
14 # %%
```

```

15 # Global Variables
16 train_directory = r"/scratch/gilbreth/dfarache/ece60146/Nikita/hw06/
    train2014"
17 val_directory = r"/scratch/gilbreth/dfarache/ece60146/Nikita/hw06/val2014"
18 train_annotations = r"/scratch/gilbreth/dfarache/ece60146/Nikita/hw06/
    annotations/instances_train2014.json"
19 val_annotations = r"/scratch/gilbreth/dfarache/ece60146/Nikita/hw06/
    annotations/instances_val2014.json"
20
21
22 categories = ["bus", "cat", "pizza"]
23 catIds = [6, 17, 59]
24 categories = {"bus": 6, "cat": 17, "pizza": 59}
25 inverse_categories = {6: "bus", 17: "cat", 59: "pizza"}
26
27
28 columns = ["id", "path_to_image", "bus", "cat", "pizza"]
29 new_image_size = 256
30 resize_image = tvf.Compose([tvf.Resize((new_image_size, new_image_size))])
31
32
33 train = False
34 coco = COCO(train_annotations) if train else COCO(val_annotations)
35 path_to_dir = "/scratch/gilbreth/dfarache/ece60146/Nikita/hw06/Train" if
    train else "/scratch/gilbreth/dfarache/ece60146/Nikita/hw06/Val"
36
37 # %%
38 def get_all_image_ids():
39     all_image_ids = []
40
41     # Get all combinations of the categories
42     for idx in range(1, len(categories) + 1):
43         combination = [list(i) for i in itertools.combinations(categories,
            idx)]
44
45         for class_labels in combination:
46             category_ids = coco.getCatIds(catNms=class_labels)
47             image_ids = coco.getImgIds(catIds=category_ids)
48
49             for image_id in image_ids:
50                 if(image_id not in all_image_ids):
51                     all_image_ids.append(image_id)
52
53     print(f"Total number of images for all the combinations of the
    categories: {len(all_image_ids)}")
54     return all_image_ids
55
56 all_image_ids = get_all_image_ids()
57
58 # %%
59 def resize_bbox(bbox, original_width, original_height):
60     # Resize bbox
61     annotation_width, annotation_height = bbox[2], bbox[3]
62     annotation_left_x, annotation_left_y = bbox[0], bbox[1]

```

```

63
64     x_scale = new_image_size / original_width
65     y_scale = new_image_size / original_height
66
67     new_bbox_width = x_scale * annotation_width
68     new_bbox_height = y_scale * annotation_height
69     new_x1 = x_scale * annotation_left_x
70     new_y1 = y_scale * annotation_left_y
71
72     new_x2 = new_x1 + new_bbox_width
73     new_y2 = new_y1 + new_bbox_height
74
75     # Add everything to the list to add to the df
76     ann_bbox = [new_x1, new_y1, new_x2, new_y2]
77     return ann_bbox
78
79 # %%
80 def resize_image_bbox(image, path_to_image, Bboxes):
81     directory = "/scratch/gilbreth/dfarache/ece60146/Nikita/hw06/train2014"
82     " if train else "/scratch/gilbreth/dfarache/ece60146/Nikita/hw06/
83     val2014"
84     new_bboxes = {"bus": [], "cat": [], "pizza": []}
85
86     try:
87         image_filename = image["file_name"]
88         original_width, original_height = image["width"], image["height"]
89
90         original_image = resize_image(read_image(os.path.join(directory,
91         image_filename)))
92         original_image = original_image.repeat(3,1,1) if original_image.
93         size()[0] == 1 else original_image # Check if image has three channels
94         original_image = tvf.functional.to_pil_image(original_image).
95         convert("RGB") # convert to RGB image
96
97         for label, bboxes in Bboxes.items():
98             for idx in range(len(bboxes)):
99                 if(bboxes[idx] != -1):
100                     resized_bbox = resize_bbox(bboxes[idx], original_width,
101                     original_height)
102
103                     if(label == "bus"):
104                         new_bboxes["bus"].append(resized_bbox)
105                         new_bboxes["cat"].append(-1)
106                         new_bboxes["pizza"].append(-1)
107
108                     elif(label == "cat"):
109                         new_bboxes["cat"].append(resized_bbox)
110                         new_bboxes["bus"].append(-1)
111                         new_bboxes["pizza"].append(-1)
112
113                     elif(label == "pizza"):
114                         new_bboxes["pizza"].append(resized_bbox)
115                         new_bboxes["cat"].append(-1)
116                         new_bboxes["bus"].append(-1)

```

```

111         original_image.save(path_to_image)
112         return new_bboxes
113
114
115     except Exception as e:
116         print(f"Exception: {e}")
117         return False
118
119 # %%
120 dataset = {}
121 for image_id in tqdm(all_image_ids):
122     image_info = {"Id": None,
123                  "Path to Image": None,
124                  "Bbox": {"bus": [], "cat": [], "pizza": []}}
125
126     image = coco.loadImgs(ids=image_id)[0]
127     annotation_ids = coco.getAnnIds(imgIds=image["id"], iscrowd=False)
128     annotations = coco.loadAnns(annotation_ids)
129
130     for annotation in annotations:
131         if(annotation["category_id"] in catIds):
132             if(annotation["area"] > (64 * 64)):
133                 path_to_image = os.path.join(path_to_dir, image["file_name
134     ])
135
136                 image_info["Id"] = image_id
137                 image_info["Path to Image"] = path_to_image
138
139                 label_name = inverse_categories[annotation['category_id']]
140                 if(label_name == "bus"):
141                     image_info['Bbox']['bus'].append(annotation['bbox'])
142                     image_info['Bbox']['cat'].append(-1)
143                     image_info['Bbox']['pizza'].append(-1)
144                 elif(label_name == "cat"):
145                     image_info['Bbox']['cat'].append(annotation['bbox'])
146                     image_info['Bbox']['bus'].append(-1)
147                     image_info['Bbox']['pizza'].append(-1)
148                 elif(label_name == "pizza"):
149                     image_info['Bbox']['pizza'].append(annotation['bbox'])
150                     image_info['Bbox']['cat'].append(-1)
151                     image_info['Bbox']['bus'].append(-1)
152
153             if(image_info["Path to Image"] != None):
154                 resized_bboxes = resize_image_bbox(image, image_info["Path to
155     Image"], image_info["Bbox"])
156                 if(resized_bboxes):
157                     image_info["Bbox"] = resized_bboxes
158                     dataset[image["file_name"]] = image_info
159
160 print(f'Total Images downloaded {len(dataset)}')
161
162 # %%
163 def create_dataframe(dataset):
164     df_filename = "train_data.csv" if train else "test_data.csv"
165     df = pd.DataFrame(columns=columns)

```

```

163
164     ids, paths_to_image, bus_bboxes, cat_bboxes, pizza_bboxes = [], [], [],
165     [], []
166     for image_filename, image_info in dataset.items():
167         # image_info = {id:str, path_to_image:str, bboxes:dict}
168         ids.append(image_info["Id"])
169         paths_to_image.append(image_info["Path to Image"])
170         bus_bboxes.append(image_info["Bbox"]["bus"])
171         cat_bboxes.append(image_info["Bbox"]["cat"])
172         pizza_bboxes.append(image_info["Bbox"]["pizza"])
173
174     assert len(bus_bboxes) == len(cat_bboxes) == len(pizza_bboxes) == len(ids
175     ) == len(paths_to_image), f"Lengths of data not matched: {len(bus_bboxes
176     )}, {len(cat_bboxes)}, {len(pizza_bboxes)}, {len(ids)}, {len(
177     paths_to_image)}"
178     df["id"] = ids
179     df["path_to_image"] = paths_to_image
180     df["bus"] = list(bus_bboxes)
181     df["cat"] = list(cat_bboxes)
182     df["pizza"] = list(pizza_bboxes)
183     df.to_csv(df_filename)
184
185 create_dataframe(dataset)
186
187 # %%

```

Listing 1: COCO Downloader

### 2.1.2 Inputs Sample

The following image is a sample of the training data used for the YOLO network

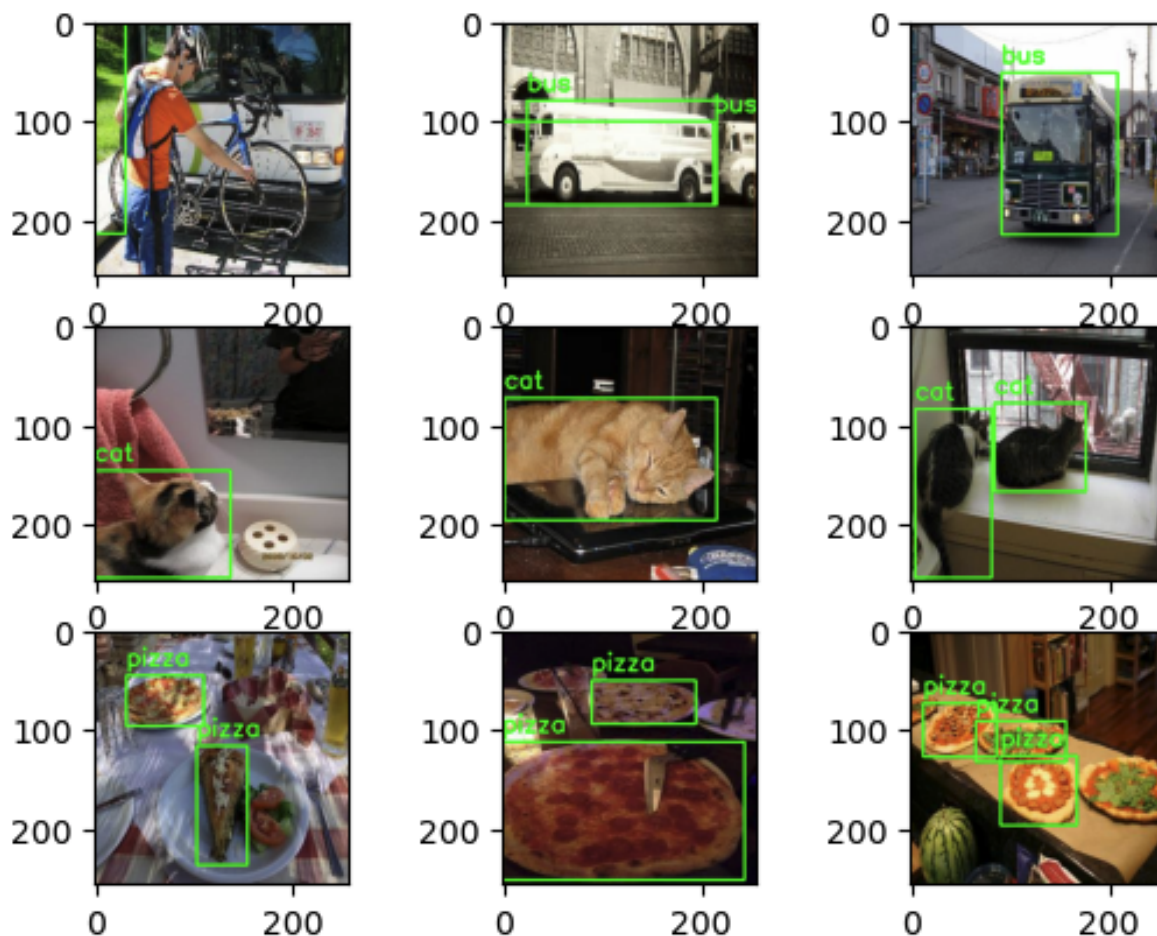


Figure 1: Inputs

## 2.2 Dataloader

### 2.2.1 Creation of the YOLO Tensor

The YOLO tensor is created through the following steps

1. Iterate through all the foreground objects present in the image downloaded from the COCO dataset
2. Overlay the image with a  $6 \times 6$  grid where each cell has a dimension of  $42 \times 42$  pixels
3. Compute the coordinates for the center of the object in the image

4. Obtain the cell coordinate that corresponds to the object found in the image
5. Assign to a variable, the height and width of the bounding box
6. Calculate the x-coordinate displacement  $\delta x$  and the y-coordinate displacement  $\delta y$  from the center of the bounding box and the center of the cell
7. Calculate the aspect ratio (bounding box height divided by the bounding box width). For this assignment, the aspect ratios used are 1/5, 1/3, 1/1, 3/1, and 5/1
8. Create a vector, called the yolo vector, of size 8 to store the indicator of an object presence,  $\delta x$ ,  $\delta y$ , the height of the bounding box, the width of the bounding box, and the one hot encoding of the class label
9. Append the yolo vector to a yolo tensor by the cell index and anchor box index
10. Finally create an augmented version of the yolo tensor to throw all the probability mass into an extra ninth element of the yolo vector if no object is present

```

1 def create_yolo_tensor(bboxes, labels, num_images_in_batch=1):
2     # Inspired by Professor Kak's
3     run_code_for_training_multi_instance_detection function
4     height_center_bb = torch.zeros(num_images_in_batch, 1).float()
5     width_center_bb = torch.zeros(num_images_in_batch, 1).float()
6     object_bb_height = torch.zeros(num_images_in_batch, 1).float()
7     object_bb_width = torch.zeros(num_images_in_batch, 1).float()
8     yolo_tensor = torch.zeros(num_yolo_cells, num_anchor_boxes,
9                               yolo_vector_size)
10
11     # jdx represents the index of the object in the foreground of the
12     image
13     for jdx in range(max_num_objects):
14         object_center_y_coord = (bboxes[jdx, 1].float() + bboxes[jdx, 3].
15         float()) / 2.0 # Just check the batch and iterate through each bbox for
16         y1 and y2
17         height_center_bb = object_center_y_coord.int()
18
19         object_center_x_coord = (bboxes[jdx, 0].float() + bboxes[jdx, 2].
20         float()) / 2.0 # Just check the batch and iterate through each bbox for
21         x1 and x2
22         width_center_bb = object_center_x_coord.int()
23
24         object_bb_height = bboxes[jdx, 3] - bboxes[jdx, 1] # Height of the
25         bounding box
26         object_bb_width = bboxes[jdx, 2] - bboxes[jdx, 0] # Width of the
27         bounding box
28
29         if(object_bb_height < 4.0 or object_bb_width < 4.0):
30             # Disregard bounding boxes that have a height or width of less
31             than an arbitrary number 4
32             continue

```

```

24     # Get the cell row and column index that corresponds to the center
    of the bounding box
25     cell_row_idx = torch.clamp((height_center_bb / yolo_interval).int
    (), max=num_cells_image_height - 1)
26     cell_col_idx = torch.clamp((width_center_bb / yolo_interval).int()
    , max=num_cells_image_width - 1)
27
28     ### Get the height of the bounding box divided by the actual
    height of the cell
29     bh = object_bb_height.float() / yolo_interval
30     bw = object_bb_width.float() / yolo_interval
31
32     ### Calculate del_x and del_y
33     # del_x is the x_coordinate displacement from the center of the
    bbox and the center of the cell
34     # del_y is the y_coordinate displacement from the center of the
    bbox and the center of the cell
35
36     # Calculate the center of the cell (i,j) coordinates
37     cell_center_i = cell_row_idx * yolo_interval + float(yolo_interval
    ) / 2.0
38     cell_center_j = cell_col_idx * yolo_interval + float(yolo_interval
    ) / 2.0
39
40     # Compute del_x and del_y
41     del_x = (object_center_x_coord.float() - cell_center_j.float()) /
    yolo_interval
42     del_y = (object_center_y_coord.float() - cell_center_i.float()) /
    yolo_interval
43
44     ### Get the class label
45     class_label_of_object = int(labels[jdx].item())
46     if(class_label_of_object == 13):
47         # Disregard labels with class label 13
48         continue
49
50     ### Get Aspect Ratio
51     aspect_ratio = object_bb_height.float() / object_bb_width.float()
52     anchor_box_idx = 0
53     if aspect_ratio <= 0.2:
54         anchor_box_idx = 0
55         ## (45)
56     if 0.2 < aspect_ratio <= 0.5:
57         anchor_box_idx = 1
58         ## (46)
59     if 0.5 < aspect_ratio <= 1.5:
60         anchor_box_idx = 2
61         ## (47)
62     if 1.5 < aspect_ratio <= 4.0:
63         anchor_box_idx = 3
64         ## (48)
65     if aspect_ratio > 4.0:
66         anchor_box_idx = 4
67
68     ### Create the yolo vector
69     yolo_vector = torch.FloatTensor([0, del_x.item(), del_y.item(), bh
    .item(), bw.item(), 0, 0, 0])
70     yolo_vector[0] = 1 # Object exists
71     yolo_vector[5 + class_label_of_object] = 1

```



```

63
64     ### Assign to yolo tensor
65     yolo_cell_index = cell_row_idx.item() * num_cells_image_width +
cell_col_idx.item()
66     yolo_tensor[yolo_cell_index, anchor_box_idx] = yolo_vector # 1
x36x5x8
67
68     ## Create an augmented yolo tensor where if no object is present,
throw all the prob mass into the extra 9th element of yolo vector
69     yolo_tensor_aug = torch.zeros(num_yolo_cells, num_anchor_boxes,
yolo_vector_size + 1).float()
70     yolo_tensor_aug[:, :, -1] = yolo_tensor
71
72     # If no object is present, throw all the prob mass into the extra 9th
element of the yolo vector
73     for icx in range(num_yolo_cells):
74         for iax in range(num_anchor_boxes):
75             if(yolo_tensor_aug[icx, iax, 0] == 0):
76                 yolo_tensor_aug[icx, iax, -1] = 1
77
78     return yolo_tensor_aug

```

Listing 2: Creation of YOLO Tensor

## 2.2.2 The PyTorch Dataset Class

The dataset class returns the image tensor, the bounding boxes in that image tensor, the labels in that image tensor, and finally its corresponding yolo tensor for the dataloaders

```

1 class GenerateDataset(torch.utils.data.Dataset):
2     def __init__(self, df, transform=None):
3         super().__init__()
4         self.df = df
5         self.transform = transform
6
7     def __return_integer_encoding(self, category):
8         categories = {"bus": 0, "cat": 1, "pizza": 2}
9         return int(categories[category])
10
11    def __convert_string_to_list(self, info):
12        try:
13            return json.loads(info)
14        except Exception as e:
15            print(info)
16
17    def __return_image(self, info):
18        path_to_image = os.path.join(r"/scratch/gilbreth/dfarache/ece60146
/Nikita/hw06", info["path_to_image"])
19        image = Image.open(path_to_image)
20        image = self.transform(image) if self.transform else image
21
22        return image
23
24    def __return_bbox_and_labels(self, info):

```

```

25     bus_bbox = self.__convert_string_to_list(info["bus"])
26     cat_bbox = self.__convert_string_to_list(info["cat"])
27     pizza_bbox = self.__convert_string_to_list(info["pizza"])
28     label_bbox_dict = {"bus": bus_bbox, "cat": cat_bbox, "pizza":
pizza_bbox}
29
30     assert len(bus_bbox) == len(cat_bbox) == len(pizza_bbox), f"Number
of annotations are not the same {len(bus_bbox)} != {len(cat_bbox)} !=
{len(pizza_bbox)}"
31
32     labels = torch.zeros(max_num_objects, dtype=torch.uint8) + 13 # 13
was randomly selected to differentiate from labelled objects in the
foreground
33     bboxes = torch.zeros(max_num_objects, 4, dtype=torch.uint8)
34
35     row = 0
36     for label, bbs in label_bbox_dict.items():
37         for jdx in range(len(bbs)):
38             bb = bbs[jdx]
39             if(bb != -1 and row < max_num_objects):
40                 bboxes[row] = torch.tensor(bb, dtype=torch.float)
41                 labels[row] = self.__return_integer_encoding(label)
42                 if(row < max_num_objects):
43                     row += 1
44
45     return labels, bboxes
46
47     def __len__(self):
48         return len(self.df)
49
50     def __getitem__(self, idx):
51         image_info = self.df.iloc[idx]
52         image = self.__return_image(image_info)
53         labels, bboxes = self.__return_bbox_and_labels(image_info)
54         yolo_tensor = create_yolo_tensor(bboxes, labels)
55
56     return image, bboxes, labels, yolo_tensor

```

Listing 3: Generate Dataset

```

1 def get_dataloader(path, debug=False):
2     # Constants
3     transform = tvn.Compose([tvn.ToTensor(), tvn.Normalize((0.5, 0.5, 0.5)
, (0.5, 0.5, 0.5))])
4
5     # Create PyTorch Datasets and Dataloader
6     df = pd.read_csv(path)
7     dataset = GenerateDataset(df, transform)
8
9     if(debug):
10         print(f"Length of dataset: {len(dataset)}")
11         print(f"Image size: {dataset[0][0].shape}")
12         print("Bounding Boxes")
13         print(dataset[0][1])

```

```

14     print("Labels")
15     print(dataset[0][2])
16     print("Yolo Tensor")
17     print(dataset[0][3])
18
19     dataloader = torch.utils.data.DataLoader(dataset, batch_size=
batch_size, num_workers=2, shuffle=True, drop_last=True)
20     return dataloader
21
22 # Get Trainloader
23 trainloader = get_dataloader(path=r"/scratch/gilbreth/dfarache/ece60146/
Nikita/hw06/train_data.csv", debug=True)
24 testloader = get_dataloader(path=r"/scratch/gilbreth/dfarache/ece60146/
Nikita/hw06/test_data.csv", debug=True)

```

Listing 4: Skip Connections

## 2.3 YOLO Algorithm

### 2.3.1 The Networks

The network designed is similar to the one used in homework 5 - single object detection and localization with a slight modification. The returned value from the network is going to be the predicted yolo vector which has a final of 1620 output features. The entire network has 80 layers and 55,140,564 learnable parameters.

```

1 class SkipBlock(nn.Module):
2     # Inspired by Professor Kak's SkipBlock class
3     def __init__(self, in_ch, out_ch, downsample=False, skip_connections=
True):
4         super(SkipBlock, self).__init__()
5         self.downsample = downsample
6         self.skip_connections = skip_connections
7         self.in_ch = in_ch
8         self.out_ch = out_ch
9         self.conv1 = nn.Conv2d(self.in_ch, self.out_ch, kernel_size=3,
stride=1, padding=1)
10        self.conv2 = nn.Conv2d(self.in_ch, self.out_ch, kernel_size=3,
stride=1, padding=1)
11        self.bn = nn.BatchNorm2d(self.out_ch)
12
13        if(downsample):
14            self.downsampler = nn.Conv2d(self.in_ch, self.out_ch,
kernel_size=1, stride=2)
15
16        def forward(self, x):
17            identity = x
18            out = self.conv1(x)
19            out = self.bn(out)
20            out = F.relu(out)
21            if(self.in_ch == self.out_ch):
22                out = self.conv2(out)
23                out = self.bn(out)

```

```

24         out = F.relu(out)
25         if(self.downsample):
26             out = self.downsampler(out)
27             identity = self.downsampler(identity)
28         if(self.skip_connections):
29             if(self.in_ch == self.out_ch):
30                 out = out + identity
31             else:
32                 out = torch.cat((out[:, :, self.in_ch :, :] + identity, out[:,
self.in_ch :, :, :] + identity), dim=1)
33         return out

```

Listing 5: Skip Connections

```

1 class NetForYolo(nn.Module):
2     # Inspired by Professor Kak's NetForYolo class
3     def __init__(self, skip_connections=True, depth=8):
4         super(NetForYolo, self).__init__()
5         self.skip_connections = skip_connections
6         self.depth = depth // 2
7         self.conv1 = nn.Conv2d(in_channels=3, out_channels=64, kernel_size
=3, padding=1)
8         self.conv2 = nn.Conv2d(in_channels=64, out_channels=64,
kernel_size=3, padding=1)
9         self.pool = nn.MaxPool2d(kernel_size=2, stride=2)
10        self.bn1 = nn.BatchNorm2d(num_features=64)
11        self.bn2 = nn.BatchNorm2d(num_features=128)
12        # self.bn3 = nn.BatchNorm2d(num_features=256)
13        self.skip64_arr = nn.ModuleList()
14        for idx in range(self.depth):
15            self.skip64_arr.append(SkipBlock(in_ch=64, out_ch=64,
skip_connections=self.skip_connections))
16
17        self.skip64ds = SkipBlock(in_ch=64, out_ch=64, downsample=True,
skip_connections=self.skip_connections)
18        self.skip64to128 = SkipBlock(in_ch=64, out_ch=128,
skip_connections=self.skip_connections)
19        self.skip128_arr = nn.ModuleList()
20        for idx in range(self.depth):
21            self.skip128_arr.append(SkipBlock(in_ch=128, out_ch=128,
skip_connections=self.skip_connections))
22
23        self.skip128ds = SkipBlock(in_ch=128, out_ch=128, downsample=True,
skip_connections=self.skip_connections)
24        # self.skip128to256 = SkipBlock(in_ch=128, out_ch=256,
skip_connections=self.skip_connections)
25        # self.skip256_arr = nn.ModuleList()
26        # for idx in range(self.depth):
27        #     self.skip256_arr.append(SkipBlock(in_ch=256, out_ch=256,
skip_connections=self.skip_connections))
28
29        # self.skip256ds = SkipBlock(in_ch=256, out_ch=256, downsample=
True, skip_connections=self.skip_connections)
30        self.fc_seqn = nn.Sequential(

```

```

31         nn.Linear(in_features=128*16*16, out_features=6*6*5*9),
32 #         nn.ReLU(inplace=True),
33 #         nn.Linear(in_features=4096, out_features=2048),
34 #         nn.ReLU(inplace=True),
35 #         nn.Linear(in_features=2048, out_features=6*6*5*9) # 6x6 grid
        overlaying the image, 5 anchor boxes, length 9 yolo vector
36     )
37
38     def forward(self, x):
39         x = self.pool(F.relu(self.conv1(x)))
40         x = self.pool(F.relu(self.conv2(x)))
41         for idx, skip64 in enumerate(self.skip64_arr[:self.depth//4]):
42             x = skip64(x)
43         x = self.skip64ds(x)
44         for idx, skip64 in enumerate(self.skip64_arr[:self.depth//4]):
45             x = skip64(x)
46         x = self.bn1(x)
47         x = self.skip64to128(x)
48         for idx, skip128 in enumerate(self.skip128_arr[:self.depth//4]):
49             x = skip128(x)
50         x = self.bn2(x)
51         x = self.skip128ds(x)
52         x = x.view(-1, 128*16*16)
53         x = self.fc_seqn(x)
54         return x

```

Listing 6: YOLO Network

### 2.3.2 Training Algorithm

The training algorithm consists of passing the image into the network shown above to obtain the predicted yolo tensor. We then compute the binary cross entropy loss for the object presence indicator, MSE loss for the regression elements of the yolo vectors in the tensor, and finally cross entropy loss on the predicted class labels. The training function accommodates for variable batch sizes through indexing of the yolo tensors and only checks yolo vectors that have an object.

```

1 def train(net, trainloader):
2     # Inspired by Professor Kak's
    run_code_for_training_multi_instance_detection function
3     net = net.to(device)
4     presence_average_loss = []
5     regression_average_loss = []
6     class_average_loss = []
7
8     print("Training started")
9
10    for epoch in range(1, epochs + 1):
11        print(f"Epoch: {epoch} / {epochs}")
12        running_bce_loss = 0.0
13        running_mse_loss = 0.0
14        running_ent_loss = 0.0
15

```

```

16         for batch_idx, (images_in_batch, _, _, yolo_tensors_in_batch) in
17             enumerate(trainloader):
18                 images_in_batch = images_in_batch.to(device)
19                 yolo_tensors_in_batch = yolo_tensors_in_batch.to(device)
20
21                 optimizer.zero_grad()
22                 output = net(images_in_batch)
23                 # print(output.shape)
24                 output = output.view(batch_size, num_yolo_cells,
25                                     num_anchor_boxes, yolo_vector_size+1)
26
27                 total_bce_loss = torch.tensor(0.0, requires_grad=True).float()
28                 .to(device)
29                 total_mse_loss = torch.tensor(0.0, requires_grad=True).float()
30                 .to(device)
31                 total_ent_loss = torch.tensor(0.0, requires_grad=True).float()
32                 .to(device)
33
34                 # Get Yolo Vectors where object presence is not zero
35                 sigmoid = nn.Sigmoid()
36                 object_presence = torch.nonzero(output[:, :, :, 0])
37                 bce_loss = criterion1(sigmoid(output[:, :, :, 0]),
38                                     yolo_tensors_in_batch[:, :, :, 0])
39                 total_bce_loss += bce_loss
40
41                 # Compute Regression Loss
42                 predicted_regression_vector = output[object_presence[:,0],
43                                     object_presence[:,1], object_presence[:,2], 1:5]
44                 target_regression_vector = yolo_tensors_in_batch[
45                                     object_presence[:,0], object_presence[:,1], object_presence[:,2], 1:5]
46                 mse_loss = criterion2(predicted_regression_vector,
47                                     target_regression_vector)
48                 total_mse_loss += mse_loss
49
50                 # Compute Cross Entropy Loss
51                 class_probs_vector = output[object_presence[:,0],
52                                     object_presence[:,1], object_presence[:,2], 5:]
53                 target_class_vector = torch.argmax(yolo_tensors_in_batch[
54                                     object_presence[:,0], object_presence[:,1], object_presence[:,2], 5:],
55                                     dim=1)
56                 cross_loss = criterion3(class_probs_vector,
57                                     target_class_vector)
58                 total_ent_loss += cross_loss
59
60                 total_ent_loss.backward(retain_graph=True)
61                 total_mse_loss.backward(retain_graph=True)
62                 total_bce_loss.backward(retain_graph=True)
63
64                 optimizer.step()
65
66                 running_bce_loss += total_bce_loss.item()
67                 running_mse_loss += total_mse_loss.item()
68                 running_ent_loss += total_ent_loss.item()

```

```

57         if(batch_idx % 20 == 19):
58             print(f"[Epoch: {epoch}/{epochs}, Batch Idx= {batch_idx +
59 1}] ----- Average BCE Loss: {running_bce_loss / float(20)}, Average
MSE Loss: {running_mse_loss / float(20)}, Average Cross Entropy Loss: {
running_ent_loss / float(20)}")
60
61             presence_average_loss.append(running_bce_loss / float(20))
62             regression_average_loss.append(running_mse_loss / float
(20))
63             class_average_loss.append(running_ent_loss / float(20))
64
65             torch.save(net.state_dict(), path_to_model)
66
67             running_bce_loss = 0.0
68             running_mse_loss = 0.0
69             running_ent_loss = 0.0
70
71         return presence_average_loss, regression_average_loss,
class_average_loss
72
73 presence_loss, regression_loss, class_loss = train(net, trainloader)

```

Listing 7: Training

The parameters to the training algorithm were

Parameters	Value
Epochs	12
Learning Rate	1e−6
Optimizer	Adam
Betas	(0.9, 0.99)
Criterion 1	Binary Cross Entropy
Criterion 2	Mean Squared Error
Criterion 3	Cross Entropy

The following image illustrates the binary cross entropy loss, the regression loss, and the cross entropy loss across all 12 epochs

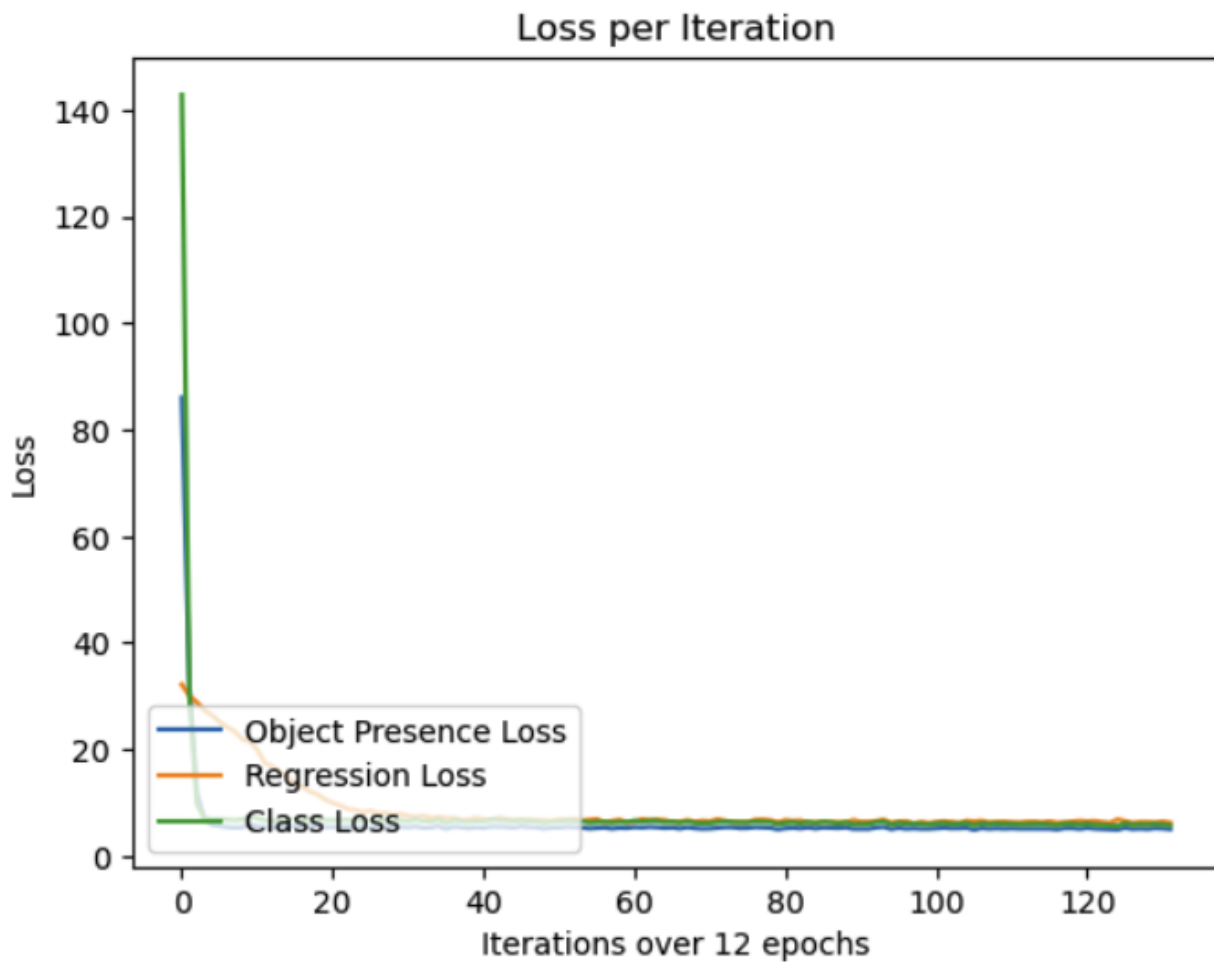


Figure 2: Training Losses

### 2.3.3 Testing Algorithm

In the testing function, we iterate across one batch, across one batch axis to obtain the predicted bounding boxes and labels for each image. A comparison of the ground truth and the predictions are illustrated in figure ???. The implementation for the testing function is shown below.

```

1 def test(net, data):
2     net.load_state_dict(torch.load(path_to_model))
3     net.eval()
4     net.to(device)
5
6     image, bboxes_in_image, labels_in_image, yolo_tensor = data
7
8     # Add an extra dimension to the image to pass it to the network
9     image_for_network = torch.unsqueeze(deepcopy(image), dim=0).to(device)
10

```



```

11     predicted_yolo_tensor = net(image_for_network)
12     predicted_yolo_tensor = predicted_yolo_tensor.view(1, num_yolo_cells,
13     num_anchor_boxes, yolo_vector_size+1)
14     yolo_tensor = torch.unsqueeze(yolo_tensor, dim=0)
15
16     predictions = {}
17     softmax = nn.Softmax(dim=0)
18
19     for icx in range(num_yolo_cells):
20         for iax in range(num_anchor_boxes):
21             predicted_yolo_vector = predicted_yolo_tensor[0, icx, iax]
22             target_label = torch.argmax(softmax(predicted_yolo_vector[5:]))
23
24             if target_label.item() != 3:
25                 predictions[(icx, iax)] = predicted_yolo_vector
26
27     # Display Image with annotations
28     image = np.asarray(tvt.ToPILImage()(image / 2 + 0.5))
29
30     for key, yolo_vector in predictions.items():
31         cell_idx, anchor_box_idx = key
32         object_presence, del_x, del_y, bh, bw, bus, cat, pizza, extra =
33         yolo_vector
34
35         predicted_label = torch.argmax(softmax(torch.tensor([bus, cat,
36         pizza, extra])))
37
38         if predicted_label == 0: label = 'bus'
39         if predicted_label == 1: label = 'cat'
40         if predicted_label == 2: label = 'pizza'
41
42         cell_col_idx = cell_idx % num_cells_image_width
43         cell_row_idx = cell_idx // num_cells_image_height
44
45         cell_center_i = cell_row_idx * yolo_interval + float(yolo_interval
46         ) / 2.0
47         cell_center_j = cell_col_idx * yolo_interval + float(yolo_interval
48         ) / 2.0
49
50         obj_center_x = del_x * yolo_interval + cell_center_j
51         obj_center_y = del_y * yolo_interval + cell_center_i
52
53         bh *= yolo_interval
54         bw *= yolo_interval
55
56         x1 = int(obj_center_x - bw / 2)
57         y1 = int(obj_center_y - bh / 2)
58
59         x2 = int(obj_center_x + bw / 2)
60         y2 = int(obj_center_y + bh / 2)
61
62         image = cv2.rectangle(image, (x1, y1), (x2, y2), color=(0, 0,
63         255), thickness=1)

```

```

58     image = cv2.putText(image, label, (x1, y1-10), cv2.
FONT_HERSHEY_SIMPLEX, 0.5, (0, 0, 255), 2)
59
60     # Ground Truth
61     for idx in range(len(bboxes_in_image)):
62         x1 = int(bboxes_in_image[idx,0])
63         y1 = int(bboxes_in_image[idx,1])
64         x2 = int(bboxes_in_image[idx,2])
65         y2 = int(bboxes_in_image[idx,3])
66
67         gt_label = int(labels_in_image[idx])
68         if(gt_label == 13):
69             continue
70         gt_label = inverse_categories[gt_label]
71
72         image = cv2.rectangle(image, (x1, y1), (x2, y2), color=(0, 255,
0), thickness=1)
73         image = cv2.putText(image, gt_label, (x1, y1-10), cv2.
FONT_HERSHEY_SIMPLEX, 0.5, (0, 0, 255), 2)
74     return image

```

Listing 8: Testing

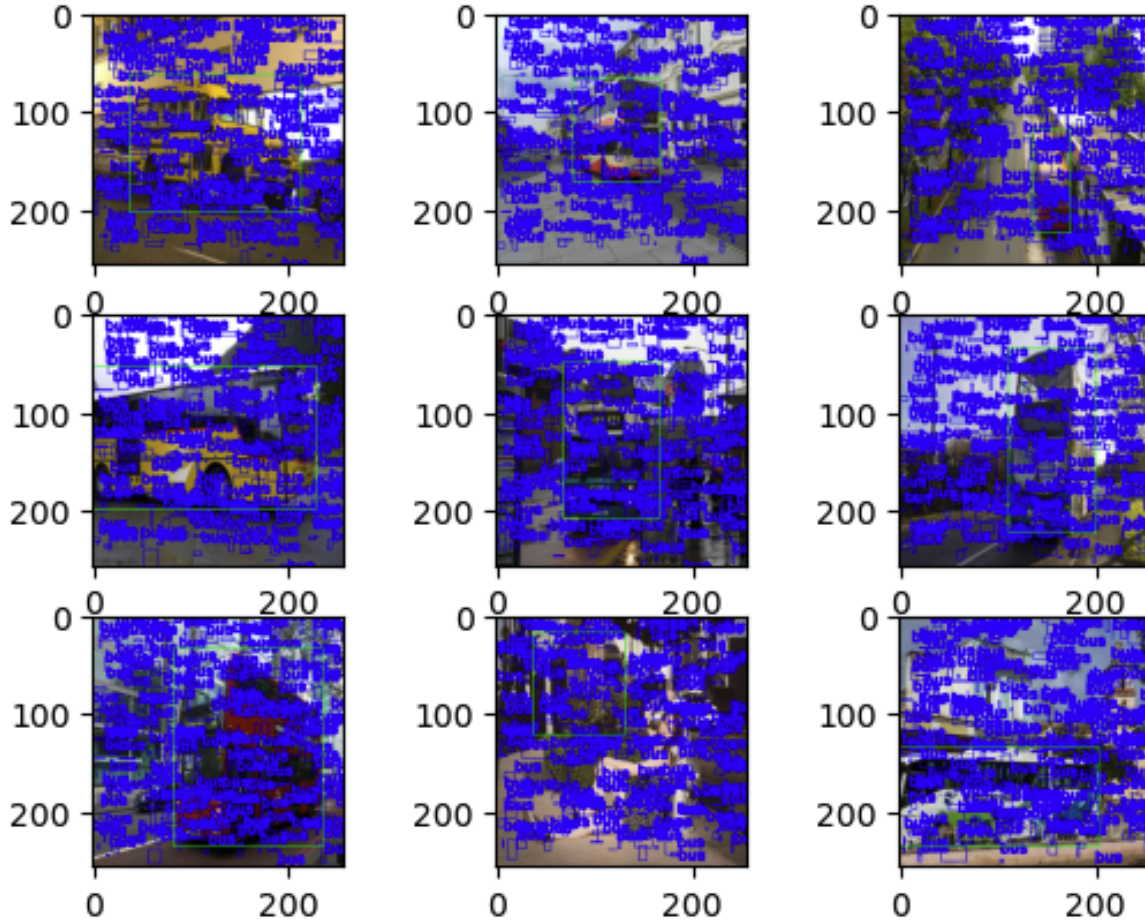


Figure 3: Inputs

### 3 Evaluation

The YOLO algorithm has a high accuracy rate for identifying the foreground objects in the images. However, with respect to the bounding boxes, the location is usually correct but the size of the bounding box is inaccurate. This could be a problem in the dataset or the way the yolo tensors are created. Maybe removing the augmented yolo tensor could improve the results.