# Video Streaming and Tracking HW2- Object Tracking Report

Author: 313581001 黃得誠

**I.     Experiment Setup**

**Data pre-process:**

Reference: https://haobin-tan.netlify.app/docs/ai/computer-vision/object-detection/coco-dataset-format/

I wrote a script called "covert2coco.py" to write the dataset information into a json file that meets the coco format. I used three columns in this assignment: "categories", " images", and "annotations".

```python
 8   def Convert2coco(label_dir, output_path):
 9       coco_data = defaultdict(list)
10       image_id = 0
11       annotation_id = 0
12
13       # Define categories
14       coco_data["categories"] = [{"id": 0, "name": "car"}]
15       labels = os.listdir(label_dir)
16       labels.sort()
17       for label_file in labels:
18           if label_file.endswith('.txt'):
19               # Extract image info
20               image_info = {
21                   "file_name": label_file.replace('.txt', '.jpg'),
22                   "height": HEIGHT,
23                   "width": WIDTH,
24                   "id": image_id
25               }
26               coco_data["images"].append(image_info)
```

```python
28               # Read annotations from the label file
29               with open(os.path.join(label_dir, label_file), 'r') as file:
30                   lines = file.readlines()
31
32               for line in lines:
33                   parts = line.strip().split()
34                   class_id, x_center, y_center, width, height = map(float, parts)
35
36                   # Convert normalized positions to absolute (pixel) positions
37                   abs_x_center = x_center * WIDTH
38                   abs_y_center = y_center * HEIGHT
39                   abs_width = width * WIDTH
40                   abs_height = height * HEIGHT
41
42                   # Convert to COCO format (x_min, y_min, width, height)
43                   x_min = abs_x_center - (abs_width / 2)
44                   y_min = abs_y_center - (abs_height / 2)
```

```
46                        annotation = {
47                            "id": annotation_id,
48                            "image_id": image_id,
49                            "category_id": int(class_id),
50                            "bbox": [x_min, y_min, abs_width, abs_height],
51                            "area": abs_width * abs_height,
52                            "iscrowd": 0
53                        }
54                        coco_data["annotations"].append(annotation)
55                        annotation_id += 1
56
57                    image_id += 1
58
59            # Write out the COCO dataset
60            with open(output_path, 'w') as json_file:
61                json.dump(coco_data, json_file)
```
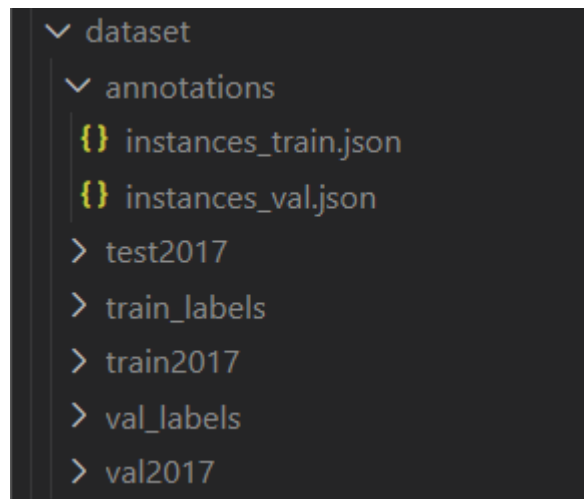
In addition, I change the name of the datasets.

Dataset:

```
∨ dataset
  ∨ annotations
    {} instances_train.json
    {} instances_val.json
  > test2017
  > train_labels
  > train2017
  > val_labels
  > val2017
```

**Hyperparameters:**

For hyperparameters, check "/exps/example/custom/yolo_s_.py".

I mostly follow the recommendation hyperparameters of the reference, which is in /exps/default/yolo_s_.py.

yolo_s_.py :

```
9   class Exp(MyExp):
10      def __init__(self):
12          self.depth = 0.33
13          self.width = 0.50
14          # self.exp_name = os.path.split(os.path.realpath(__file__))[1].split(".")[0]
15          self.exp_name = "yolox_s_SE"
16
17          # Define yourself dataset path
18          self.data_dir = "datasets/dataset"
19          self.train_ann = "instances_train.json"
20          self.val_ann = "instances_val.json"
21
22          self.num_classes = 1
23
24          self.max_epoch = 100
25          self.data_num_workers = 4
26          self.eval_interval = 1
27
28          self.basic_lr_per_img = 0.01 / 64.0
29
30          self.save_history_ckpt = False
```

Due to lack of computational resources, I use batch size=16 and train with 100 epochs. I've also tried to change the learning rate due to the change of batch size. However, it doesn't result in better results. Therefore, I finally use the recommended learning rate=0.01/64 * batch size as reference.

| | |
|---|---|
| seed | None |
| output_dir | './YOLOX_outputs' |
| print_interval | 10 |
| eval_interval | 1 |
| dataset | None |
| num_classes | 1 |
| depth | 0.33 |
| width | 0.5 |
| act | 'silu' |
| data_num_workers | 4 |
| input_size | (640, 640) |
| multiscale_range | 5 |
| data_dir | 'datasets/dataset' |
| train_ann | 'instances_train.json' |
| val_ann | 'instances_val.json' |
| test_ann | 'instances_test2017.json' |
| mosaic_prob | 1.0 |

| | |
|---|---|
| mixup_prob | 1.0 |
| hsv_prob | 1.0 |
| flip_prob | 0.5 |
| degrees | 10.0 |
| translate | 0.1 |
| mosaic_scale | (0.1, 2) |
| enable_mixup | True |
| mixup_scale | (0.5, 1.5) |
| shear | 2.0 |
| warmup_epochs | 5 |
| max_epoch | 100 |
| warmup_lr | 0 |
| min_lr_ratio | 0.05 |
| basic_lr_per_img | 0.00015625 |
| scheduler | 'yoloxwarmcos' |
| no_aug_epochs | 15 |
| ema | True |

| | |
|---|---|
| weight_decay | 0.0005 |
| momentum | 0.9 |
| save_history_ckpt | True |
| exp_name | 'yolox_s' |
| test_size | (640, 640) |
| test_conf | 0.01 |
| nmsthre | 0.65 |

II.   **Code Explanation**

I added the SE Block and Inception Module in to "yolox/models/network_blocks.py". Then add it respectively to the CSPDarkNet, which is inside the YOLOPAFPN module. YOLO_X is contained of YOLOPAFPN and YOLOXHEAD. Since the two module is used to extract feature, I decided to add them into the FPN part.

SEBlock and InceptionModule.

```python
212  class SEBlock(nn.Module):
213      def __init__(self, channel, reduction=16):
214          super(SEBlock, self).__init__()
215          self.avg_pool = nn.AdaptiveAvgPool2d(1)
216          self.fc = nn.Sequential(
217              nn.Linear(channel, channel // reduction, bias=False),
218              nn.ReLU(inplace=True),
219              nn.Linear(channel // reduction, channel, bias=False),
220              nn.Sigmoid()
221          )
222
223      def forward(self, x):
224          b, c, _, _ = x.size()
225          y = self.avg_pool(x).view(b, c)
226          y = self.fc(y).view(b, c, 1, 1)
227          return x * y.expand_as(x)
228
     You, 5 hours ago | 1 author (You)
229  class InceptionModule(nn.Module):
230      def __init__(self, in_channels, ksize=1, stride=1, act="silu"):
231          super(InceptionModule, self).__init__()
232          self.branch1x1 = BaseConv(in_channels, in_channels // 4, ksize=1, stride=stride, act=act)
233          self.branch3x3 = BaseConv(in_channels, in_channels // 4, ksize=3, stride=stride, act=act)
234          self.branch5x5 = BaseConv(in_channels, in_channels // 4, ksize=5, stride=stride, act=act)
235          self.branch_pool = BaseConv(in_channels, in_channels // 4, ksize=1, stride=stride, act=act)
236
237      def forward(self, x):
238          branch1x1 = self.branch1x1(x)
239          branch3x3 = self.branch3x3(x)
240          branch5x5 = self.branch5x5(x)
241          branch_pool = nn.functional.avg_pool2d(x, kernel_size=3, stride=1, padding=1)
242          branch_pool = self.branch_pool(branch_pool)
243          outputs = [branch1x1, branch3x3, branch5x5, branch_pool]
244          return torch.cat(outputs, 1)
```

Add SEBlock into "dark2", "dark3", "dark4", "dark5", respectively.

```python
121              # dark2
122              self.dark2 = nn.Sequential(
123                  Conv(base_channels, base_channels * 2, 3, 2, act=act),
124                  CSPLayer(
125                      base_channels * 2,
126                      base_channels * 2,
127                      n=base_depth,
128                      depthwise=depthwise,
129                      act=act,
130                  ),
131                  SEBlock(base_channels * 2),
132                  # InceptionModule(base_channels*2),
133              )
```

After checking the validation score, I chose to use SEBlock only.

**Training command:**
python tools/train.py -f exps/example/custom/yolox_s.py -d 1 -b 16 --fp16 -o
-c weights/yolox_s.pth --cache ram
The model is trained on NVIDIA GeForce RTX 2080 Ti with batch size=16 and
pretrained weight: yolox_s.pth.

**Inference command:**

python tools/demo.py image -f exps/example/custom/yolox_s.py -c
weights/best_ckpt_SE.pth --path ./datasets/dataset/test2017/ --conf 0.5 --
nms 0.45 --save_result --device gpu --fp1
Used best_ckpt_SE.pth checkpoint to inference.

I modified tools/demo.py to get the required result format.
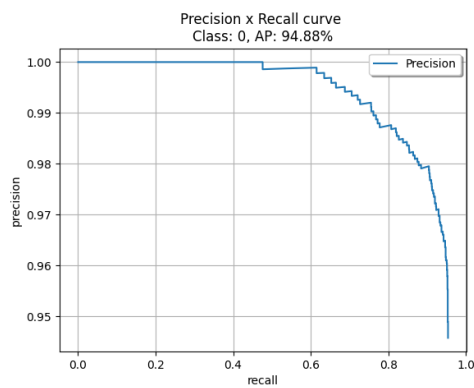Image_demo() is modified and process_outputs() is added.

```python
186  def image_demo(predictor, result_folder, path, current_time, save_result):
187      result_folder = './myresult/test'
188      if os.path.isdir(path):
189          files = get_image_list(path)
190      else:
191          files = [path]
192      files.sort()
193      for image_name in files:
194          outputs, img_info = predictor.inference(image_name)
195          # process the outputs to suitable format
196          results = process_outputs(outputs[0], img_info, predictor)
197
198          if save_result:
199              # Save the results to a text file instead of an image
200              txt_file_name = os.path.join(result_folder, os.path.splitext(os.path.basename(image_name))[0] + ".txt")
201              with open(txt_file_name, 'w') as f:
202                  for item in results:
203                      f.write("%s\n" % ' '.join(map(str, item)))
204              logger.info("Results saved to {}".format(txt_file_name))
```

```python
206  def process_outputs(outputs, img_info, predictor):
210      results = []
211      ratio = img_info["ratio"]
212      if outputs is None:
213          return results
214
215      outputs = outputs.cpu()
216
217      # Process each detection
218      for output in outputs:
219          # Extract data
220          bbox = output[0:4]  # The bounding box
221          bbox /= ratio
222          score = output[4] * output[5]  # The confidence score
223          cls = output[6]  # The class
224
225          left = bbox[0].item()
226          top = bbox[1].item()
227          right = bbox[2].item()
228          bottom = bbox[3].item()
229
230          # Add to results
231          results.append([0, score.item(), int(left), int(top), int(right), int(bottom)])
```
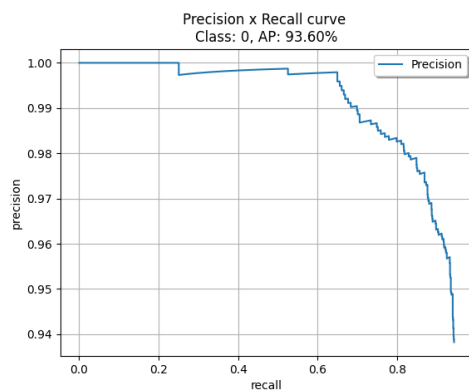
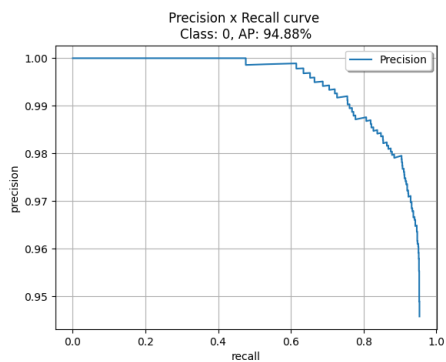## III.    Screenshot of validation results

### With SE block:



### With Inception module:



### Original:



Result: Original(Without adding extra module) > With SE module > With Inception module.

The test result is generated by the checkpoint with SE module.

**IV.  Reference**

1. https://github.com/Megvii-BaseDetection/YOLOX
2. https://github.com/rafaelpadilla/Object-Detection-Metrics
3. https://haobin-tan.netlify.app/docs/ai/computer-vision/object-detection/coco-dataset-format/

checkpoints with SE module: Code/YOLOX/weights/best_ckpt_SE.pth
environment: Code/YOLOX/requirements.txt