# HW3 Object Detection + Data Augmentation Report

R11945005 郭庭沂

1) Image captioning

    a. Compare three different pre-trained models



| model | Generated text |
|---|---|
| Salesforce/blip2-opt-2.7b | this is a picture of a starfish in a tank |
| Salesforce/blip2-opt-6.7b-coco | this is a picture of a fish tank with a starfish and rocks |
| Salesforce/blip2-flan-t5-xl | this is a picture of a starfish in an aquarium |

As the result, "Salesforce/blip2-opt-6.7b-coco" generated more details of the image than the others. Therefore, we will choose "Salesforce/blip2-opt-6.7b-coco" as the pretrained model for later problems.

    b. Design 2 templates of prompts for later generating comparison.

```
▼ 100
    image  "IMG_3123_jpeg_jpg.rf.e76afbe597d98c497fcf2ac25c3d0fc7.jpg"
  ► label  []  1 item
    height  1024
    width  768
  ► bboxes  []  5 items
    generated_text  "a fish tank with a starfish and rocks"
    prompt_w_label  "a fish tank with a starfish and rocks ,starfish, height:1024, width:768"
    prompt_w_suffix  "a fish tank with a starfish and rocks ,starfish, ocean, undersea background, HD quality, highly detailed"
```
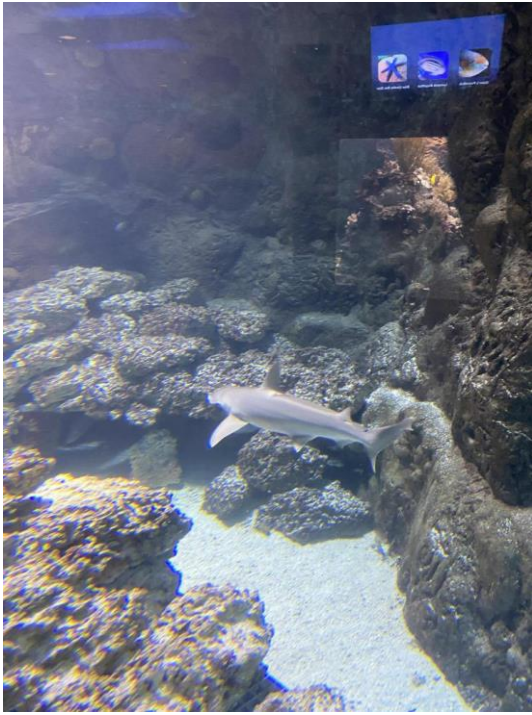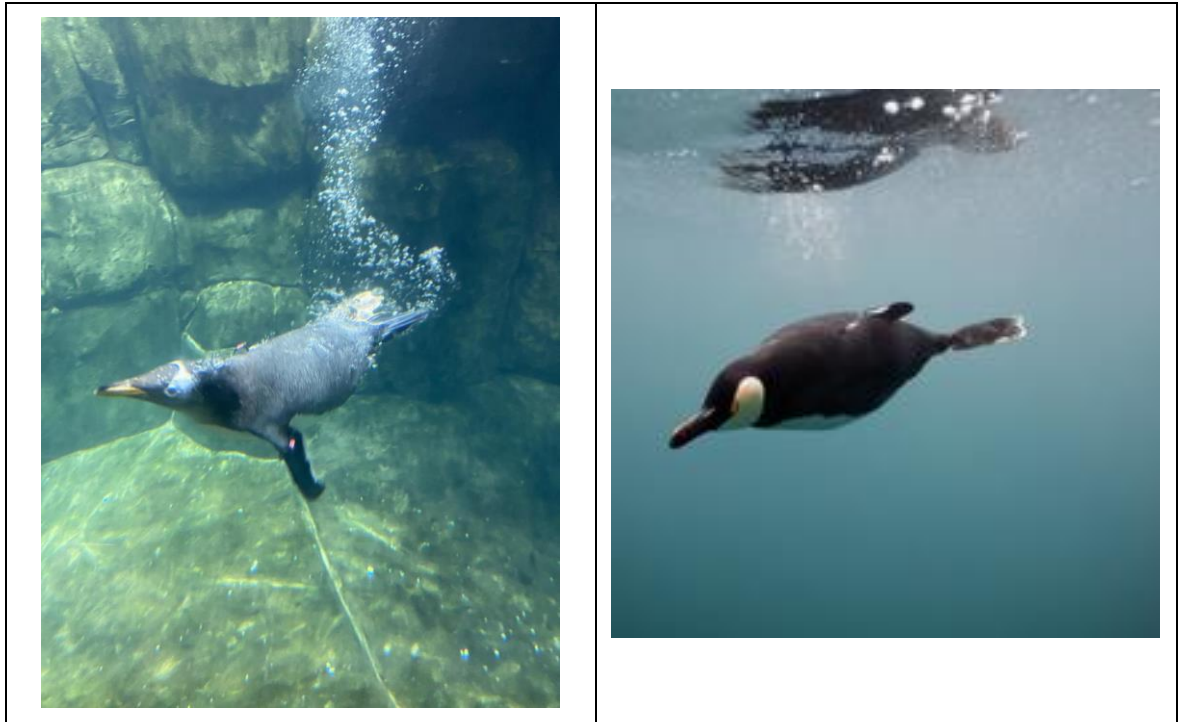
2) Text to image Generation

    a. Use two kinds of generated prompts to generate images (text only)

| Original | Prompt_w_label | Prompt_w_suffix |
|:---:|:---:|:---:|
|  |  |  |
|  |  |  |
|  |  |  |

b. Select "Prompt_w_label" prompt, and perform image grounding generation. (text + image)

| Original | "Prompt_w_label" + image |
|:---:|:---:|

3) Table of performance based on FID

| Prompt | Text Grounding | | Image Grounding |
|---|---|---|---|
| | Prompt_w_label | Prompt_w_suffix | Prompt_w_label |
| FID | 143 | 145 | 138 |

4) Table of the improvement of your detection model from HW1 after data augmentation



In order to solve the data imbalance problem, each category contains nearly 240 unique images by doing image grounding data augmentation, and each category contains nearly 200 unique images except for fish having original 240 images by doing text grounding data augmentation.

Number of unique image
after text-grounding  data augmentation

300
250
200
150
100
50
0

fish  jellyfish  penguin  puffin  shark  starfish  stingray

Number of unique image
after image-grounding data augmentation

300
250
200
150
100
50
0

fish  jellyfish  penguin  puffin  shark  starfish  stingray

We use the checkpoint which is pretrained on COCO 2017 object detection dataset to train models on three datasets, which are "Before Data Augmentation", "After text-grounding data augmentation", and "After image-grounding data augmentation" individually. The performance of each model on validation dataset is listed in table below.

| | Before Data Augmentation | After Data Augmentation (Text grounding) | After Data Augmentation (Image grounding) |
|---|---|---|---|
| AP | 0.5213 | 0.4931 | 0.5018 |

The results show that the performance of the models which are trained on augmented datasets are slightly degraded compared to non-augmented datasets.

There are two reasons we consider why the performance doesn't increase as expected. First, the original dataset contains 448 images, and the augmented dataset contains 1468, and 1676 images. It means that most images engaged in model training are generated images, therefore the performance of the model highly depends on the performance of the generated images. However, the data augmentation process includes image captioning and image generating, the quality of the generated images depend on the correctness of image captioning and the performance of the image generating model. We had found some miscaptioning text and some generated images seems unlikely to any categories. Also, the data imbalance of validation set might affect the evaluation of model performance too.
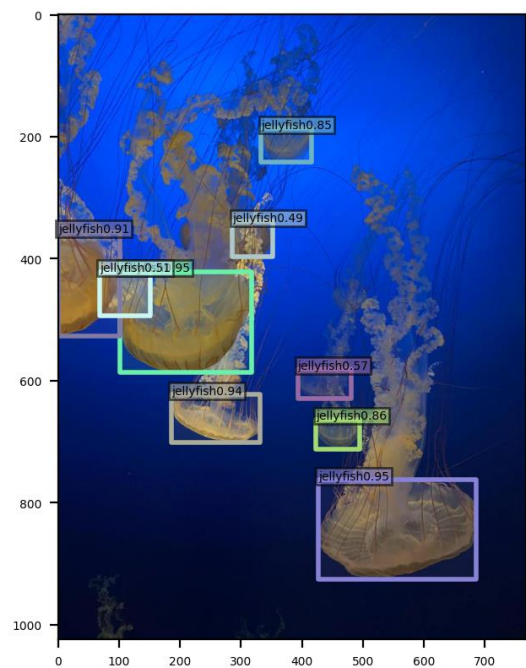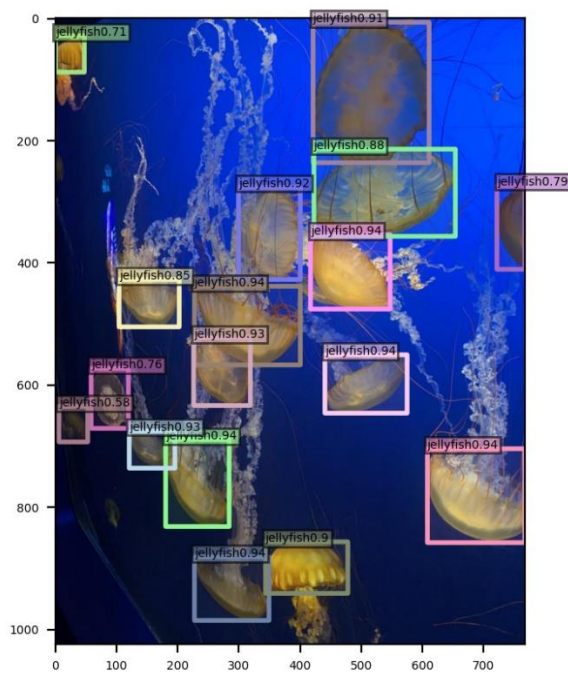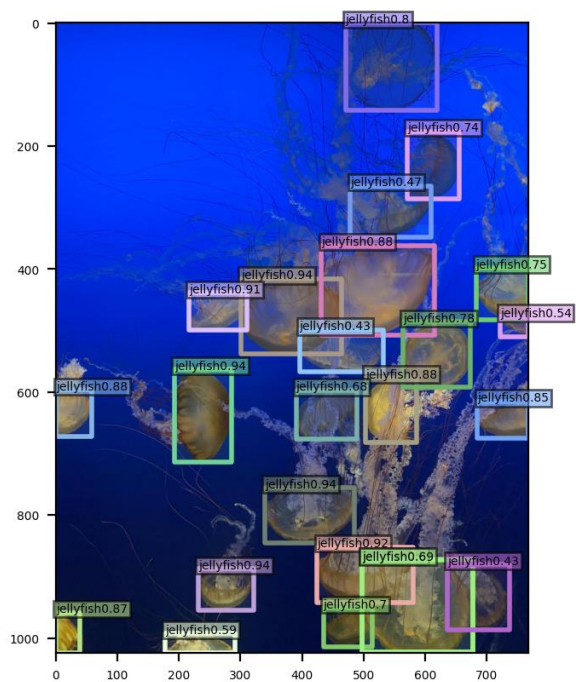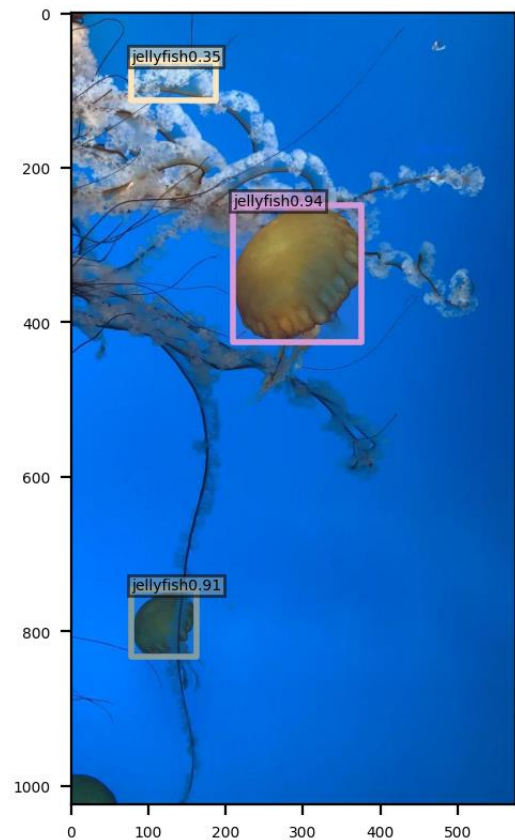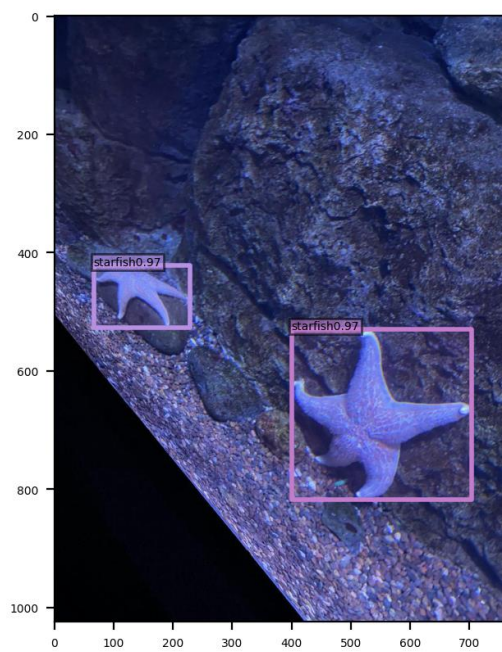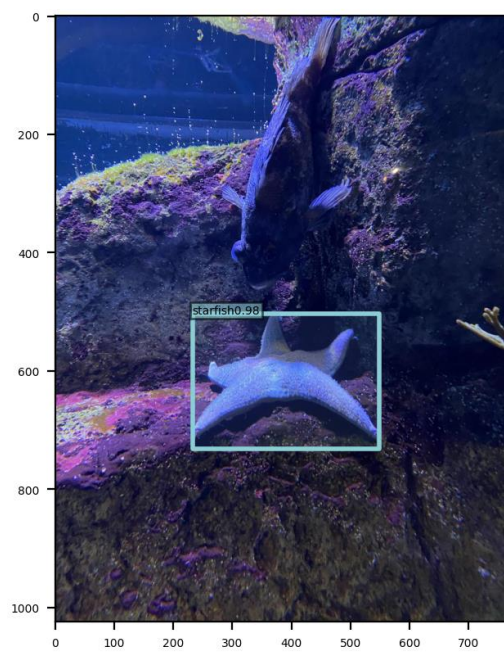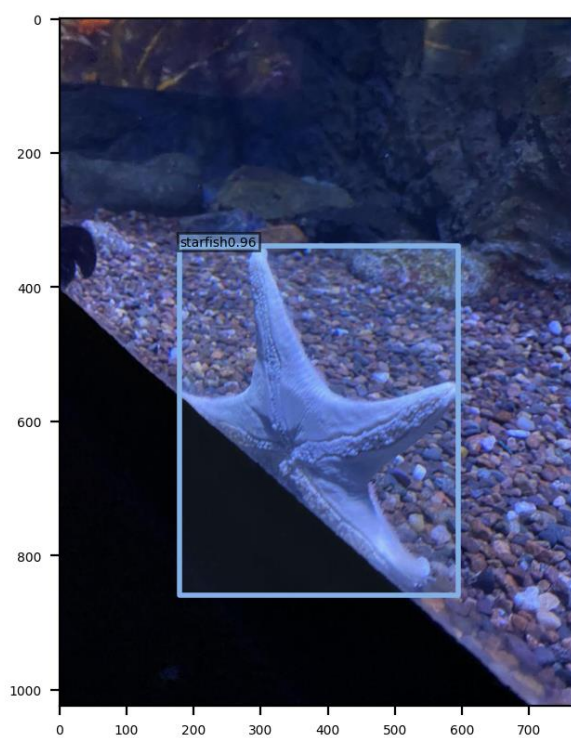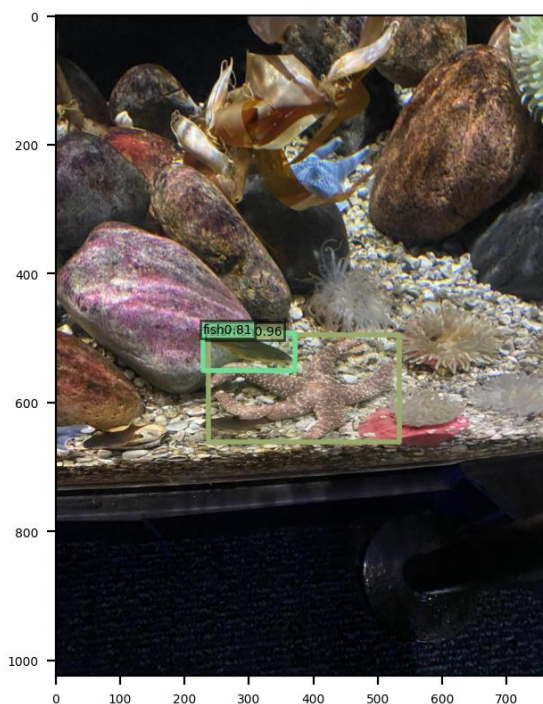
5) Visualization

    a.    Fish

b. Jellyfish

c.    Penguin

d. Puffin

e.    Shark

f.  Starfish

g.    Stingray