

# DIP Final Project – Object Detection

姓名:陳韋澎 班級:智能系統研究所 學號: 411581008

**\*\*完成程式碼與實驗結果請見:**

[https://github.com/IAII-411581008/dip\\_final\\_proj/tree/main](https://github.com/IAII-411581008/dip_final_proj/tree/main)

## 一、資料集介紹

COCO 資料集是一個大規模的影像資料集常用於物件檢測與語義分割等任務。對於從事目標檢測的學者和開發人員而言，這是一個必不可少的數據集。其具有超過 20 萬張的影像資料、80 個種類的資料共大型影像模型進行權重訓練。因此，本次實驗採用 COCO 資料集的所有物件檢測資料進行神經網路模型訓練。測試資料集將以課堂所提供之影片(person\_dog.mp4)進行模型預測。

## 二、實驗過程

### i) 作業系統環境建置

考量到實驗室與個人資源有限，並且受到桌電記憶體與效能的限制，因此最終採用 WSL2 (Windows Subsystem for Linux)而非傳統雙系統的方式進行作業系統建置。WSL2 的安裝非常簡單，僅需於 cmd 內輸入 wsl--install 並下載 Ubuntu LTS 22.04 即可。然而，在安裝 CUDA 時需要特別注意的是 NVIDIA CUDA 的驅動器僅需安裝於 Windows 端即可，WSL2 端無需另外安裝 CUDA 驅動器。考量到 CUDA 與 NVIDIA RTX 3080 顯示卡的支援度，本次實驗採用 CUDA 11.2 與 CUDNN 8.5.0 用以加速模型矩陣運算。詳細開發環境與套件版本規格，如: Python3.10, Matlab R2022B, Cmake 3.22.1, Pytorch, Scipy 等如下表所示:

圖表 1、開發環境之軟體與硬體版本與規格

硬體/軟體套件	規格/版本
GPU	NVIDIA RTX 3080 (10GB)
CPU	Intel i9-12900K@3.20GHz
DRAM	64GB
CUDA	11.2
CUDNN	8.5.0
Python	3.10
Matlab	R2022b
Cmake	3.22.1
Matplotlib	3.3.0
Numpy	1.22.2
Opencv-python	4.6.0

Pillow	7.1.2
Pyyaml	5.3.1
Requests	2.23.0
Scipy	1.4.1
Torch	1.8.0
Torchvision	0.9.0
Tqdm	4.64.0

## ii) 建立程式集開發架構

完整程式集開發架構如下圖所示：

```

1  [ultralytics]/
2  docker/
3  docs/
4  examples/
5  path/
6  runs/
7      detect/
8          yolov8x_final/ # final prediction results
9              frame2video_pre_frames/ # predicted frames
10                  0.png
11                  1.png
12                  ...
13
14              labels/ # predicted annotations
15                  frame2video_pre_1.txt
16                  frame2video_pre_1.txt
17                  ...
18              frame2video_pre.avi # raw predicted video
19
20      frame2video_pre/ # convert into video
21          frame2video_pre.avi
22      img2video/
23          person_dog.avi # predicted video w/o preprocessing
24          person_dog_yolov8x_final # final predicted video w/ preprocessing
25      plot_count_label/ # convert into frames
26          xx.png
27          ...
28      TA_result/ # convert into frames
29          xx.png
30          ...
31      video2frame_pre/ # convert into frames
32          xx.png
33          ...
34      video2images/ # convert into frames
35          xx.png
36          ...
37
38      tests/ # for env testing
39      ultralytics/ # build model utils
40
41      person_dog.mp4 # raw testing video
42      count_num_classes.py # overall process of labels' counting
43      img2video.py # convert into video
44      plot_count_label.py # count each label
45      video2img.py # convert into frames
46      video_preprocessing.mlx # exp img preprocessing
47      yolov8x.pt # YOLOv8X pretrained weights
48      requirements.txt # install required toolkits

```

圖 1、程式集開發架構圖

## iii) 測試影片前處理

我們觀察到在原始影像的光線明亮度降低(如圖 2 所示)，且部分人物與背景的對比對低。於圖 2(右)原始直方圖所示，原始影像整體的 intensities 集中分佈於 150 以下，這使得神經網路模型難以辨識出遠處樹林旁對比度與亮度降低的人類(如圖 3(左)所示)。因此，我們對於原始影像進行

了一系列處理。首先，我們觀察到影像中的人正好都穿著深色系服裝，並且正好都露出其小腿，深色服裝與皮膚的顏色在影像強度上具有一定的差異性，由此我們希望透過 Laplacian Filtering 的方式去強化其邊緣部分，使得這些特征得以被模型提取。我們透過一 4-neighborhood 的 Laplacian 濾波器，設定其 filter coefficients 為 $[0, 1, 0; 1, -4, 1; 0, 1, 0]$ 進行空間濾波，提取其高頻資訊，在將這些高頻資訊加入原始影像中，以此強化物體邊緣。然而，考量到影像中大部分的高頻成分將集中在草地的部分，為了避免過度強化草的邊緣進而惡化模型的檢測結果，在反復進行試驗後，我們最終將高頻遮罩的 coefficient 設定為 0.3。接著，從原始影像直方圖可觀察到影像整體強度較低，而 Power-law Transformation 的方式恰好適用於擴大低強度範圍並且壓縮高強度範圍。因此，我們利用了 Gamma Correction 技術，將其 Gamma Coefficient 設定為 0.8 以此改善影像中部分對比度較低的現象。透過對比圖 3 影像處理前與處理後的結果，可觀察到樹林旁的二人與其背景的對比度明顯提升了。最後，我們對於影片中的 181 個 Frames 進行影像前處理，以此強化模型的特征提取結果。影像前處理的完整程式碼如圖 6 所示。

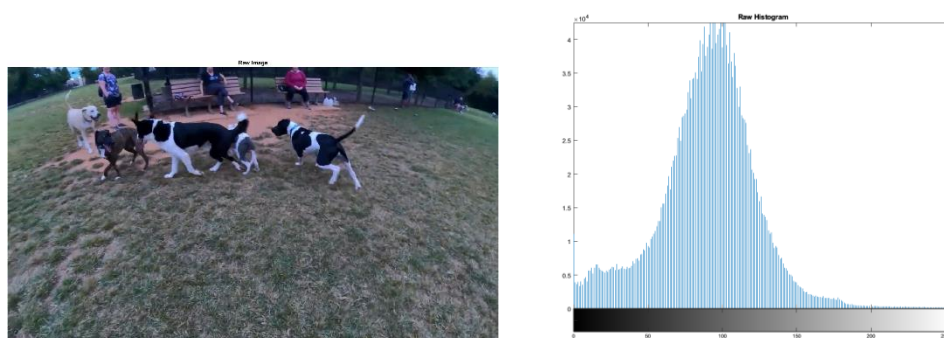


圖 2、原始第 168 Frame 的影像(左)與其原始 Histogram。



圖 3、影像處理前(左) v.s. 影像處理後(右)的明亮與對比度



圖 4、Laplacian Filtering 前(左) v.s. Laplacian 邊緣偵測(右)之影像

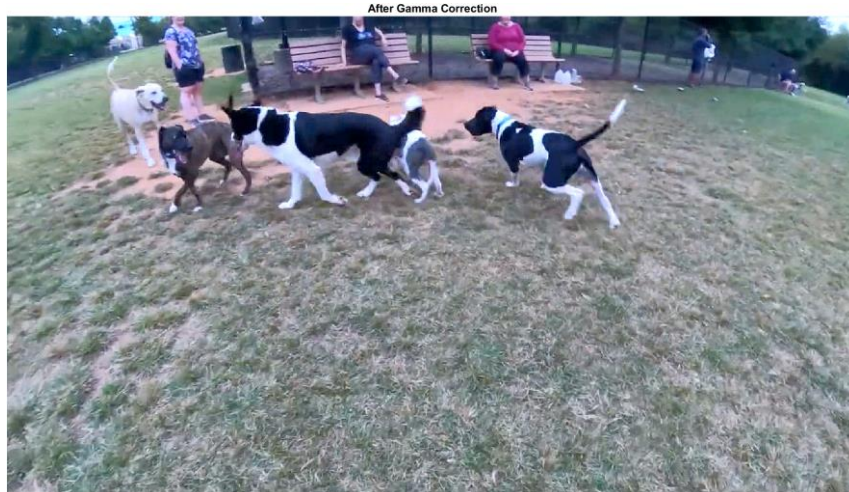


圖 5、第 168 Frame 影像前處理後的結果

```

1 % Specify the path to the video file
2 videoFilePath = 'person_dog.mp4';
3
4 % Create a VideoReader object
5 videoReader = VideoReader(videoFilePath);
6
7 % Create a folder to save the frames
8 outputFolder = 'runs/video2frame_pre';
9 if ~exist(outputFolder, 'dir')
10     mkdir(outputFolder);
11 end
12
13 % Loop through each frame and save it as an image
14 frameCount = 1;
15 while hasFrame(videoReader)
16     % Read the frame
17     frame = readFrame(videoReader);
18
19     % Convert the input image to double datatype
20     double_img = im2double(frame);
21
22     % Use the Laplacian filter to highlight edges
23     L_filter = fspecial('laplacian', 0.3); % Defining Laplacian filter
24     filtering = imfilter(double_img, L_filter, 'replicate'); % Applying Laplacian filter to image
25     filtered_img = double_img - filtering;
26
27     imshow(filtering*3)
28
29     % Selecting a scaling factor for high boost filtering
30     S = 0.3;
31
32     % Applying high boost filtering to Laplacian filtered image
33     sharpen_img = double_img + S * filtered_img;
34
35     % Applying Gamma Correction
36     gammaImg = imadjust(sharpen_img,[],[],0.8);
37
38     % imshow(gammaImg)
39     % title('After Gamma Correction')
40
41     % imhist(gammaImg)
42     % title('Histogram After Gamma Correction')
43
44     % Create the filename for the current frame
45     filename = sprintf('frame_%04d.png', frameCount);
46
47     % Save the frame as an image
48     imwrite(J, fullfile(outputFolder, filename));
49
50     % Increment the frame count
51     frameCount = frameCount + 1;
52 end
53
54 disp('Frames extracted successfully.');
```

圖 6、影像前處理程式碼

#### iv) 測試影片後處理

完成測試樣本預測後，我們需在每一幀均附上各個種類成功檢測數量，具體流程如下：

1. 將原始影片 person\_dog.mp4 以 30 FPS 分解成 181 個 Frames
2. 透過 YOLOv8X 的預測結果，計算各個 Frame's Annotation 中物品種類的數量
3. 將學號、Person 及 Dog 數量以綠色標籤 plot 與各個 Frame 左上角
4. 最後，將這些標籤好的 Frames 重新轉檔回影片.avi 檔

```

plot_num_classes.py > plot_label_counts_on_image
You, 3 hours ago | 1 author (You)
1 import cv2
2 import os
3
4 ''' video2images '''
5 def video_to_images(video_path, output_folder):
6     # Create output folder if it doesn't exist
7     if not os.path.exists(output_folder):
8         os.makedirs(output_folder)
9
10    # Open the video file
11    cap = cv2.VideoCapture(video_path)
12
13    # Get the frames per second (fps) and frame width/height
14    fps = cap.get(cv2.CAP_PROP_FPS)
15    width = int(cap.get(cv2.CAP_PROP_FRAME_WIDTH))
16    height = int(cap.get(cv2.CAP_PROP_FRAME_HEIGHT))
17
18    # Get the total number of frames in the video
19    total_frames = int(cap.get(cv2.CAP_PROP_FRAME_COUNT))
20
21    print(f"FPS: {fps}, Width: {width}, Height: {height}, Total Frames: {total_frames}")
22
23    # Loop through each frame and save it as an image
24    for frame_num in range(total_frames):
25        ret, frame = cap.read()
26
27        if not ret:
28            break
29
30        # Save the frame as an image
31        image_path = os.path.join(output_folder, f"person_dog_{frame_num:04d}.png")
32        cv2.imwrite(image_path, frame)
33
34        # Display progress
35        if frame_num % 100 == 0:
36            print(f"Processing frame {frame_num}/{total_frames}")
37
38    # Release the video capture object
39    cap.release()
40
41    print(f"Frames saved to {output_folder}")
42
43    ''' count each label '''
44    def count_labels(annotation_file):
45        # Read the annotation file and extract labels
46        with open(annotation_file, 'r') as file:
47            lines = [line.strip().split() for line in file]
48
49        # Extract labels from the first column
50        labels = [line[0] for line in lines]
51
52        # Count occurrences of labels "0" and "16"
53        label_counts = {"Person": labels.count("0"), "Dog": labels.count("16")}
54
55        # Print the counts
56        for label, count in label_counts.items():
57            print(f"{label}: {count}")
58

```

```

59 import matplotlib.pyplot as plt
60 from PIL import Image, ImageDraw, ImageFont
61
62 def plot_label_counts_on_image(image_path, annotation_file, frameIdx):
63     # Read the annotation file and extract labels
64     with open(annotation_file, 'r') as file:
65         lines = [line.strip().split() for line in file]
66
67     # Extract labels from the first column
68     labels = [line[0] for line in lines]
69
70     # Count occurrences of labels "0" and "16"
71     label_counts = {"Person": labels.count("0"), "Dog": labels.count("16")}
72
73     # Create an image with the counts displayed in the top left corner
74     original_image = Image.open(image_path)
75
76     # Create an ImageDraw object to add text to the image
77     draw = ImageDraw.Draw(original_image)
78
79     # Specify the font and size for the text (increased font size)
80     font_size = 30
81     font = ImageFont.truetype("arial.ttf", font_size)
82
83     # Add text indicating the count of each label in the top left corner with bold, red text and line breaks
84     text = ""
85     draw.text((10, 10), f"ID: {411581008} \n", font=font, fill=(0, 255, 0))
86     text += "\n"
87     for label, count in label_counts.items():
88         # Render the text twice for a bold effect
89         if label == "Person":
90             draw.text((10, 60), f"{label}: {count}", font=font, fill=(0, 255, 0))
91         else:
92             draw.text((10, 110), f"{label}: {count}", font=font, fill=(0, 255, 0))
93         # draw.text((11, 11), f"{label}: {count}", font=font, fill=(255, 0, 0))
94
95     final_image_path = os.path.join("runs\\plot_count_label", f"person_dog_{frameIdx:04d}.png")
96     original_image.save(final_image_path)
97
98     # Display the final image
99     # original_image.show()
100
101     ''' images2video '''
102     def frames_to_video(frames_folder, output_video_path, fps):
103         # Get the list of image files in the frames folder
104         image_files = [f for f in os.listdir(frames_folder) if f.endswith(('png', 'jpg', 'jpeg'))]
105
106         if not image_files:
107             print("No image files found in the specified folder.")
108             return
109
110         # Sort the image files based on their numerical names
111         image_files.sort(key=lambda x: int(''.join(filter(str.isdigit, x))))
112
113         # Get the dimensions of the first image
114         first_image = cv2.imread(os.path.join(frames_folder, image_files[0]))
115         height, width, _ = first_image.shape
116
117         # Create a VideoWriter object
118         fourcc = cv2.VideoWriter_fourcc(*'mp4v')
119         video_writer = cv2.VideoWriter(output_video_path, fourcc, fps, (width, height))
120
121         # Iterate through image files and write frames to the video
122         for image_file in image_files:
123             frame = cv2.imread(os.path.join(frames_folder, image_file))
124             video_writer.write(frame)
125             # print(image_file)
126         # Release the video writer
127         video_writer.release()
128
129     print(f"Video successfully created: {output_video_path}")

```



```

131 if __name__ == "__main__":
132     ## video2frames
133     # Specify the path to the input video and output folder
134     video_path = r"runs\detect\yolov8x_final\frame2video_pre.avi"
135     output_folder = r"runs\video2images"
136
137     # Convert the video to images
138     video_to_images(video_path, output_folder)
139
140     ## count each label
141     annotationDir = r"E:\research_proj\dip_final_proj\ultralytics\runs\detect\yolov8x_final\labels"
142     imgDir = r"E:\research_proj\dip_final_proj\ultralytics\runs\video2images"
143
144     annotationFiles = os.listdir(annotationDir)
145     # Sort the image files based on their numerical names
146     annotationFiles.sort(key=lambda x: int(''.join(filter(str.isdigit, x))))
147
148     imgFileNames = os.listdir(imgDir)
149
150     for fileId in range(0, len(annotationFiles)):
151         annotationFile = os.path.join(annotationDir, annotationFiles[fileId])
152         imgFile = os.path.join(imgDir, imgFileNames[fileId])
153         # Specify the path to your annotation file
154         # annotation_file = r"E:\research_proj\dip_final_proj\ultralytics\runs\detect\predict11\labels\perso
155
156         # Count and print the occurrences of labels "0" and "16"
157         count_labels(annotationFile)
158
159         # Plot the label counts on the image
160         plot_label_counts_on_image(imgFile, annotationFile, fileId)
161
162     ## convert back to video
163     frames_folder = r"runs\plot_count_label"
164     output_video_path = r"runs\imgs2video\person_dog_yolov8x_final.mp4"
165     fps = 30 # Specify the desired frames per second (fps) for the output video
166
167     # Replace "path/to/frames" with the path to your folder containing frames
168     frames_to_video(frames_folder, output_video_path, fps)

```

圖 7、影像後處理程式碼

### 三、預訓練模型與超參數設定

本次實驗採用 YOLOv8X 進行物件檢測。YOLO 系列的模型常用於物件檢測任務，其具有非常好檢測結果的同時，整體運算速度也相較於 RCNN-based 模型來的快。因此，我們最終採用最新版 YOLO 進行物件檢測。YOLOv8 自行開發了 C2f 模塊，並替換了 Darknet-53 模型中的 C3 模塊。損失函數的部分，其採用了 TaskAlignedAssigner 的方式進行標籤分配與 Distribution Focal Loss DIOU NMS Loss 與 CIOU Loss 進行權重調整。具體模型架構可參考圖 8。

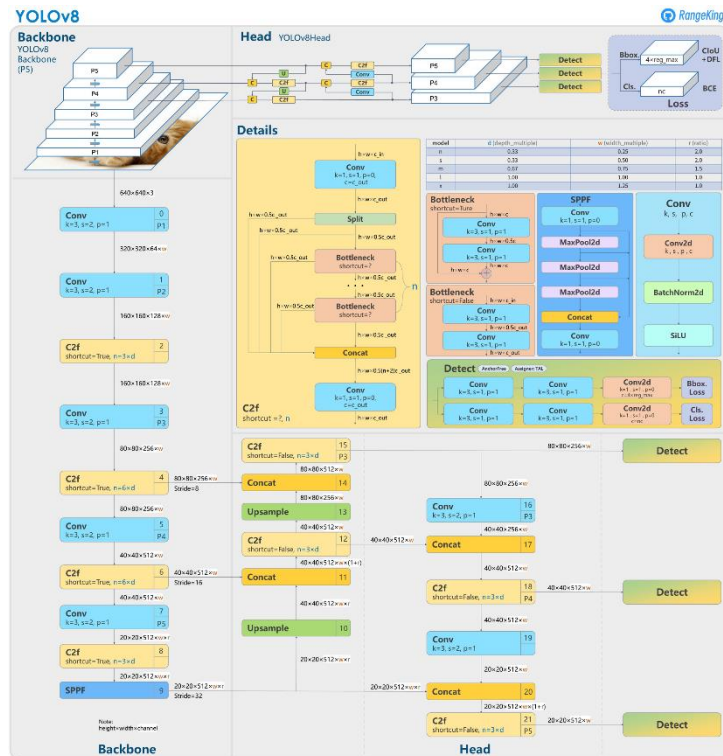


圖 8、Yolov8 架構圖 (Source: [Roboflow](#))

圖表 2 訓練超參數設定

超參數種類	參數值
Model Weights	yolo8x.pt
Confidence Score Threshold	0.2
IoU Threshold	0.5
Classes	0, 1, 16
Image Size	640
Max_det	300
Vid_stride	False
Agnostic_nms	False
Augment	False



#### 四、實驗結果與討論

實驗結果的部分，我們劃分成兩個部分進行討論：(1)對比使用影像處理後對於模型檢測結果的影響；(2)與助教的實驗結果進行比較。圖 9 為 YOLOv8X 在相同模型權重與超參數設定下，對於影像處理前與處理後的檢測結果。透過對比圖 9 的十張圖片可發現，右上角樹林旁的二人在經過 Laplacian Sharpening 與 Power-law Transformation 後更容易被 YOLOv8X 檢測出來。其主要的原因是這兩個人的腿部特征受到 High-pass Filtering 後，其邊緣部分可明顯與背景劃分開來，並且在 Gamma Correction 的幫助下擴大 low-level intensities range，使得這兩個人深色衣服的特征與陰暗背景的對比度得以進一步提升，進而促使這兩個人的特征得以成功被模型提取。另外，我們亦可對比影像處理前與處理後人物的 Confidence Score 差異性。圖八的結果所示，經影像處理後，人物整體的 Confidence Scores 明顯高於原始影像的結果，尤其是在右上角樹林下的兩個人的 Confidence Scores。然而，影像處理後亦可能帶來負面的影響，如：狗狗整體的 Confidence Scores 在進行影像處理後的表現略低於原始影像，這可能是因為 Laplacian 銳利化了草地的特征，進而導致不必要的“雜訊”被放大，影像模型的檢測結果。因此，我們在設計 Laplacian Filtering 時特意限制了其 Mask's Coefficient，避免過量的高頻資訊惡化模型的檢測結果。但透過對比左半部分與右半部分的結果，影像處理對於狗狗的影響並未非常明顯，整體 Confidence Scores 依舊保持在 80% 以上。最後，圖 9(i)、(j)的結果所示，經影像處理後 YOLOv8X 得以成功檢測出所有人物與狗狗，特別是右上角的二人与狗狗。

第二部分我們同樣的任意挑選個五張圖片與助教的實驗結果進行對比。透過對比左半部分與右半部分的結果可以發現，我們的研究方法對於人物的檢測結果更加的敏銳，如圖 10(d)的結果所示(中間凳子上的二人)，本研究方式得僅需部分人物特征(i.e., 小腿與鞋子的特征)便可成功檢測出相關人物。其可能原因有幾種：(1)本次實驗使用的影像處理技術得以有效地強化膚色、鞋子與背景的對比度，使得這些特征得以成功被模型提取；(2)本次實驗所採用的模型為最新版 YOLO 模型，模型整體權重與提取技術由於過往 YOLO 模型；(3)本次實驗所設定的 Confidence Score Threshold “可能”低於助教的預設值，因此這些物件標籤得以在較低的 Confidence Score 下被檢測出來。透過圖 10(f)、(g)的實驗結果並結合圖 9 的實驗結果得以證實上述(1)、(2)的推理，我們在採用了 YOLOv8X 模型後亦可檢測出樹林下的人物(右上角第二個人)，並且在進行影像強化後，模型最終得以檢測出右上角最右邊的人物及狗狗並同時提升了人物整體的 Confidence Score。因此，本研究方式整體對於右上角樹林下的二人及狗狗檢測結果較為精準。另外，本研究方式對於重疊物體的檢測率較由於助教的方式，如圖 10(g)、(h)的實驗結果所示(中間偏左的三隻狗狗與一個人)，助教的研究結果並未檢測出後面的三隻狗狗及人物，而本研究方式不僅能檢測出人物前的兩只狗狗，對於人物後的狗狗亦能精確地成功檢測。最後，圖 10(i)、(j)所示為 5 人 6 狗的影像，助教之研究結果漏檢了中間重疊的一隻狗及右上角遠處的狗狗及人物，而本研究方式得以成功檢測出所有人物與狗狗。由此可見，在使用 YOLOv8X 並進行影像強化後得以提升模型整體的檢測結果。



(a)



(b)



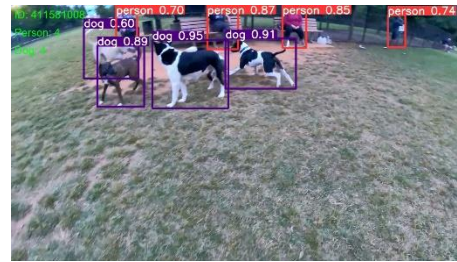
(c)



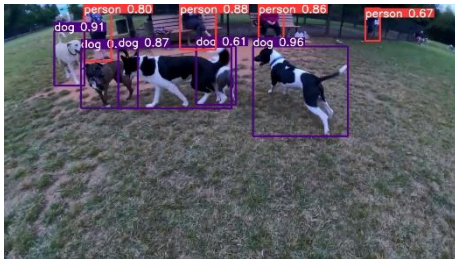
(d)



(e)



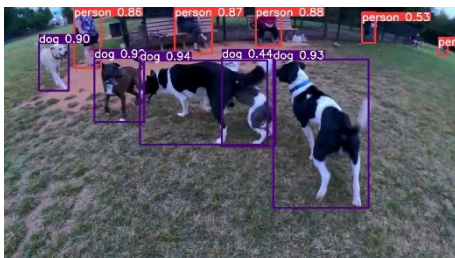
(f)



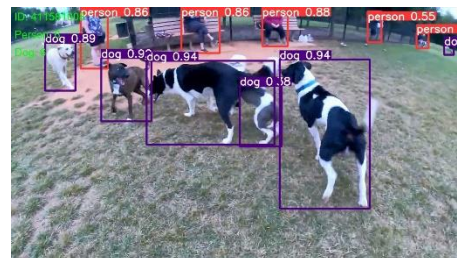
(g)



(h)



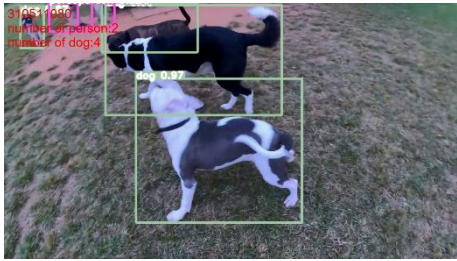
(i)



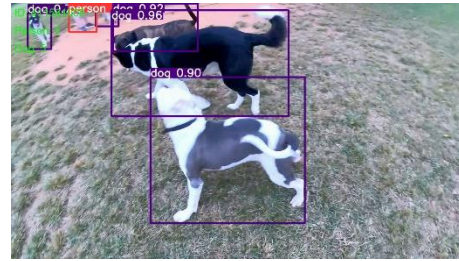
(j)

圖 9、YOLOv8X 對於 person\_dog.mp4 物件檢測結果。(a)、(c)、(e)、(g)、(i)為採用原始影像的檢測結果; (b)、(d)、(f)、(h)、(j)為採用影像處理後的檢測結果

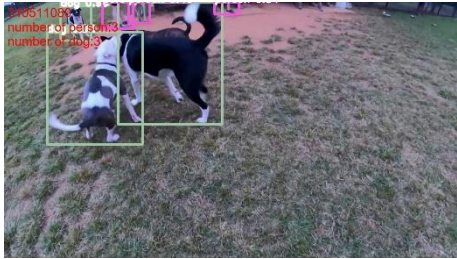




(a)



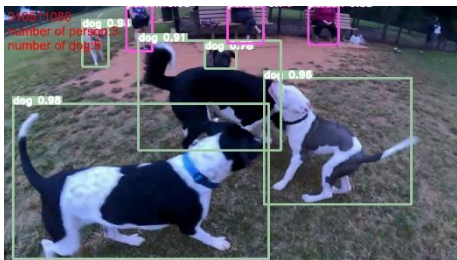
(b)



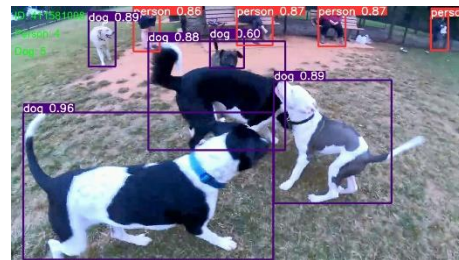
(c)



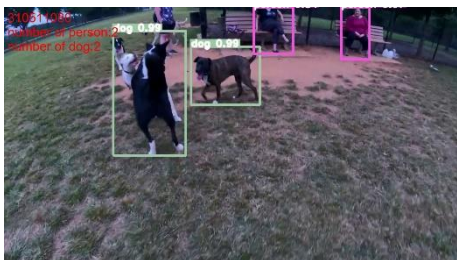
(d)



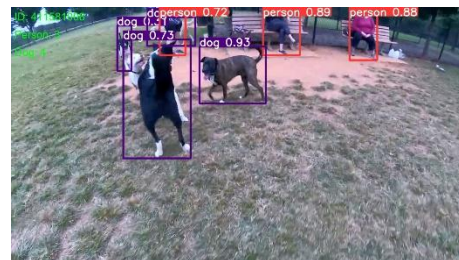
(e)



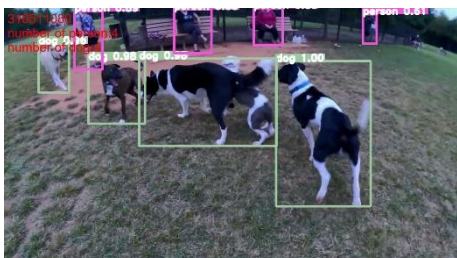
(f)



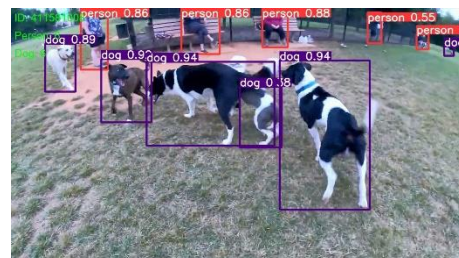
(g)



(h)



(i)



(j)

圖 10、person\_dog.mp4 物件檢測結果。(a)、(c)、(e)、(g)、(i)為助教的檢測結果；(b)、(d)、(f)、(h)、(j)為本次實驗的檢測結果

