



Object Detection (Person Tracking)

CV - C Isaac Newton

18 October 2024

Meet the Team

Dziand Dafi G.

Faster RCNN



Dziand Dafi Ginandjar

Zoe Muhammad

YOLO

Aufa Biahdillah

Faster RCNN



Aufa Biahdillah

Ryan Aprianto

YOLO



Ryan Aprianto

M Rizal Fauzi

YOLO



mochamad rizal fauzi

Background & Problem Statement

- As the population grow, industries encounter an exponential escalation in their production and services, which requires safety, efficiency, and automation.
- In quality control, for example, checking if their products are excellent for sale is essential to preserving the industry's reputation. Meanwhile, the tremendous number of products may cause an error in this quality control. Thus, by implementing object detection, industries can easily overcome this issue.
- The person-tracking model, specifically, can assist in maintaining security, such as detecting or even tracking suspicious persons in CCTV.
- There are still many applications of the person-tracking model. One of which is in the autonomous vehicle.



Objectives & Scope

- The objectives of this project are:
 - Conducting object detection for person tracking experiments applying two models: Faster RCNN and YOLO.
 - From those experiments, the best algorithm is concluded.
- For study purposes, this project is limited to:
 - The validation dataset of the COCO dataset.
 - The detection focuses on the person only; hence, the images used are filtered to those that contain at least one person.
 - The number of images is merely 500 images.

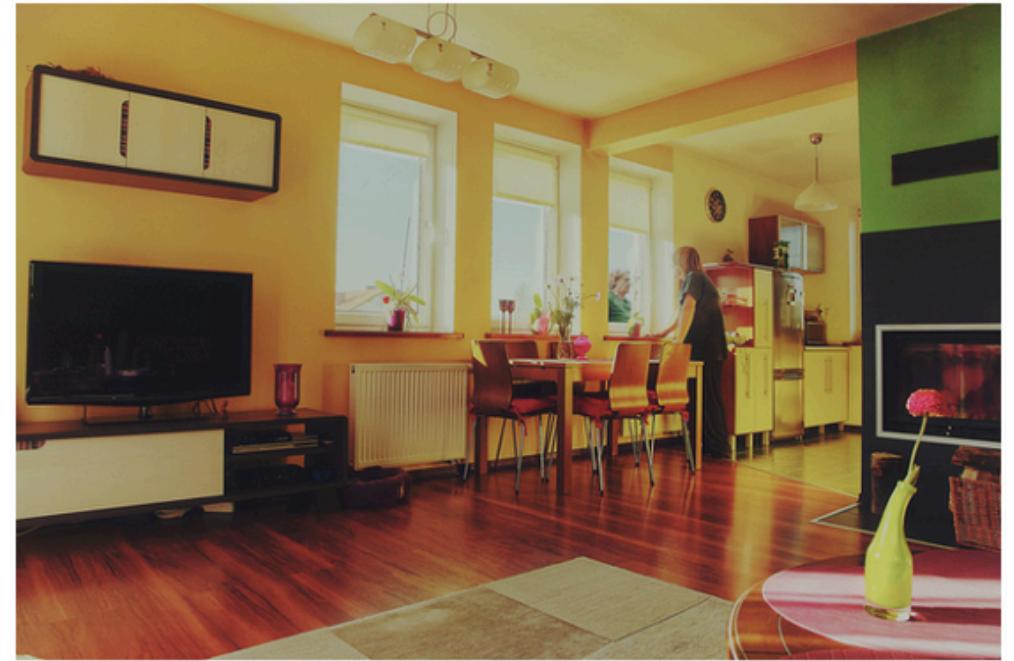


Data Collection and Exploration

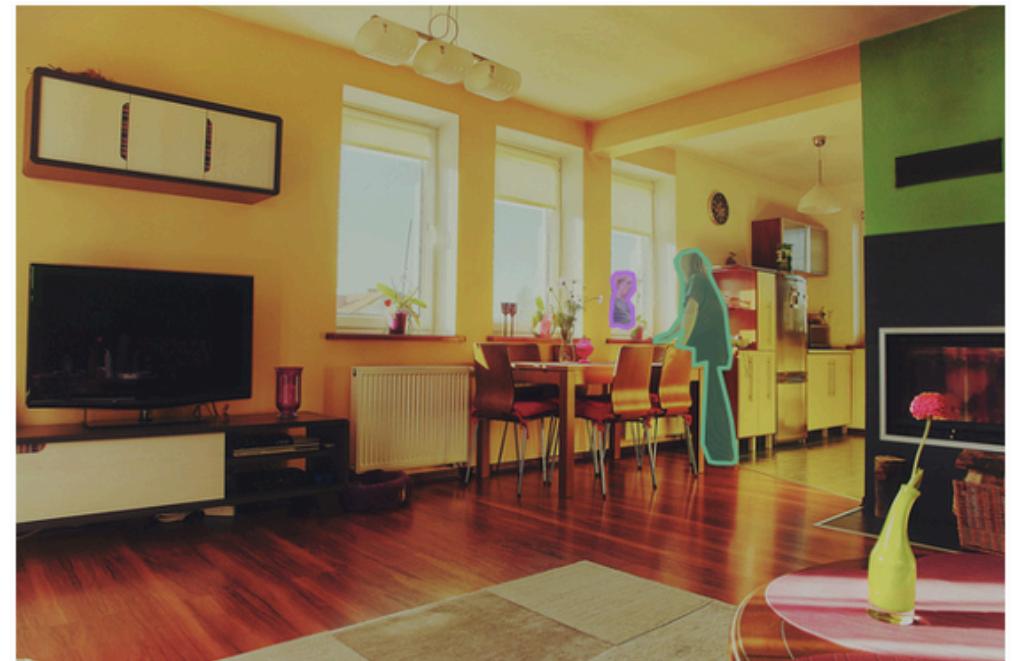
We use the COCO API, which supports us on the:

- Filtering the images into those that contain at least one person.
- Getting the information on the pictures includes the image ID, file name, and the URL used to download the image.
- Obtaining the annotation data for the image, which consists of segmentation, area, crowd, bounding box, category ID, etc., for every object in the picture.
- Filtering the annotation to person only.
- Visualizing the original and masked image for interpretation.

Original Image



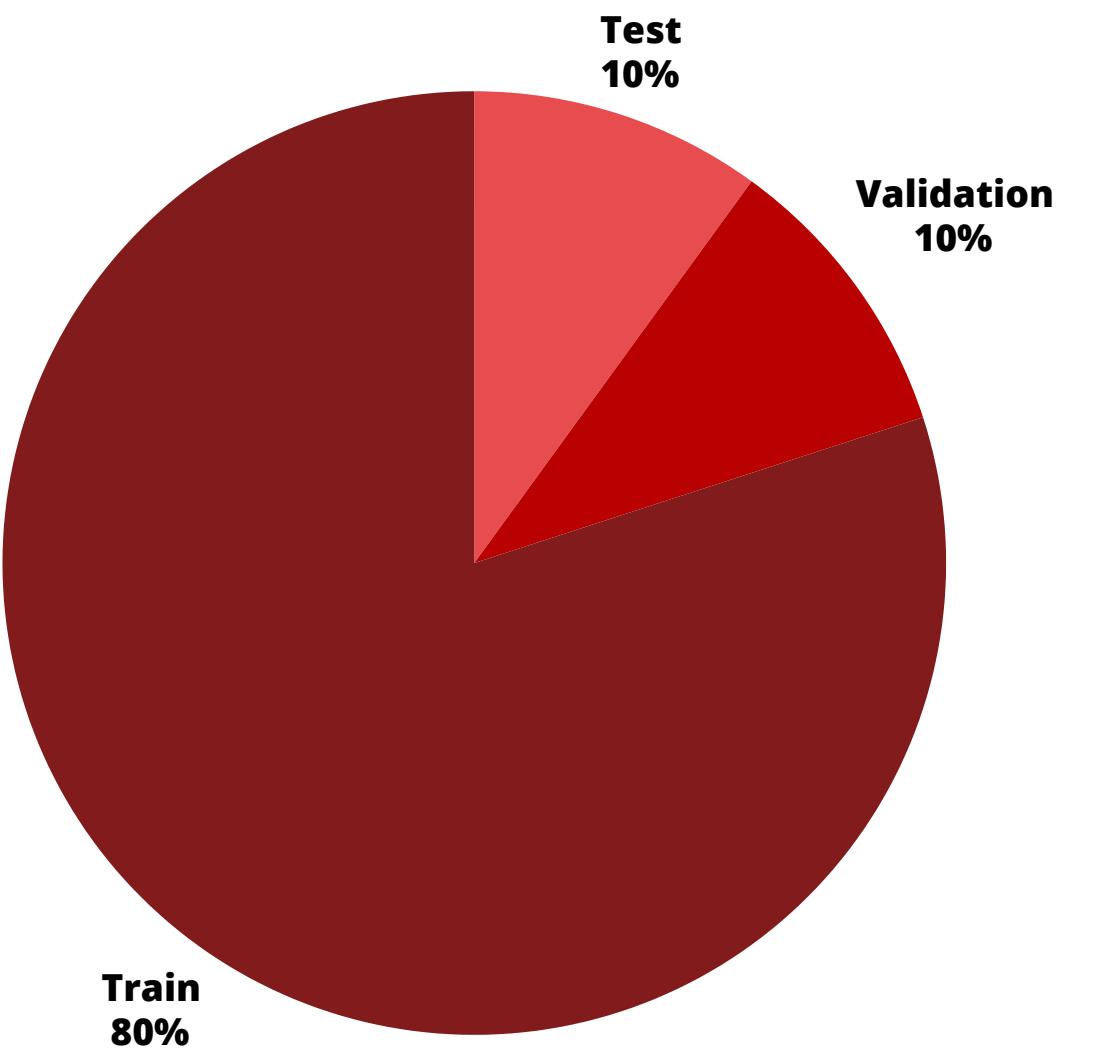
Masked Image



Data Splitting

We split the data into train, validation, and test datasets with the proportions as follows:

- 400 images for train dataset
- 50 images for validation and test dataset



Faster RCNN

Data Preparation

- We create the Cocodataset class to prepare the dataset.
- In this class, we use COCO API again to convert the annotation file into a binary mask, which is then stored in a folder.
- We also convert the bounding box format from the annotation file into XYXY format.
- Finally, we store all of the required data, which are boxes, masks, labels, image_id, area, and is_crowd, into a target variable and return it all together with the image file that has been converted to tensor and transformed if the transform function exists.

Original Image



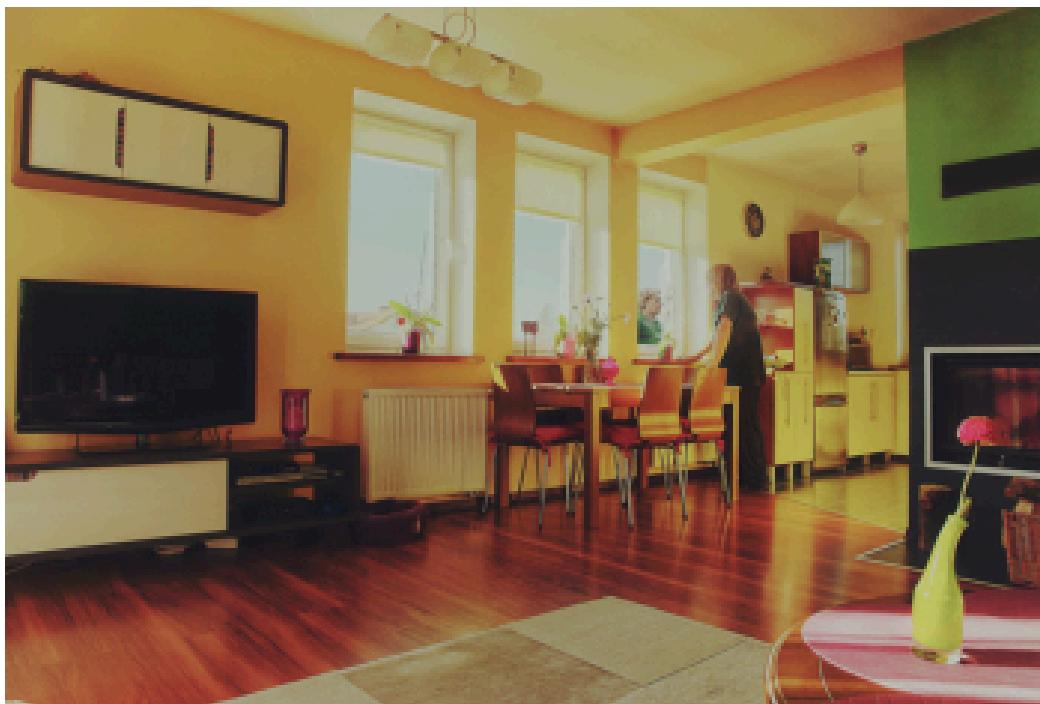
Binary Masked Image



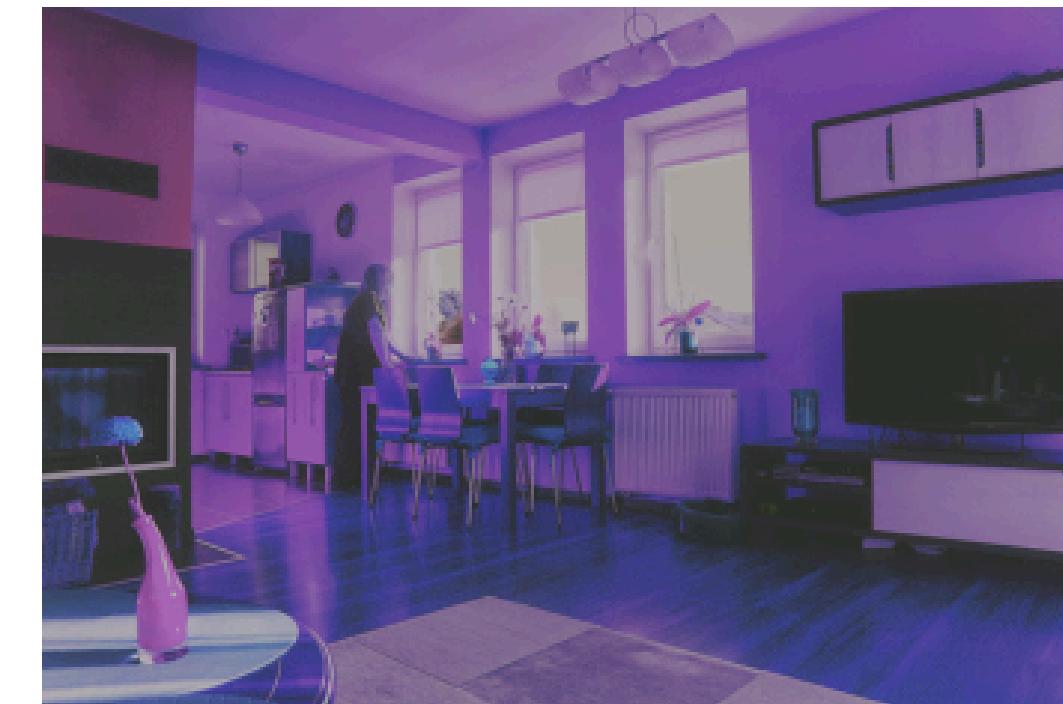
Data Transformation

- We perform the transformation to increase the images' variety, allowing the model to learn from different perspectives of the same image and improve its ability to generalize to unseen data.
- We apply the random photometric distortion, which adjusts the image's brightness, contrast, saturation, and hue.
- We also implement the horizontal flip, which flips the image horizontally.

Original



Transformed



Data Loader

We create the data loader for the datasets as:

- In the training dataset, we use `batch_size = 2` to group two samples into one batch to increase time productivity during training, while in test and validation data, we use `batch_size = 1`.
- In the training dataset, we use `shuffle = True` to shuffle the dataset at the beginning of each epoch, ensuring that the order in which samples are fed into the model changes every time. Shuffling the data is essential to prevent the model from learning the data order, which could introduce bias. In the test and validation dataset, we use `shuffle = False`.
- We set `num_workers = 2` so that two parallel processes will load data in the background while the model trains so that it can speed up the data loading.
- `collate_fn` to collate the data into a batch, that is, how individual samples from the dataset are combined into a batch.



Model Construction

We set the following for the model and its training process.



Faster RCNN

Applying the pre-trained model.



Epoch

We set the number of epoch as 50



Optimizer

SGD with $Lr = 0.005$,
 $momentum = 0.9$, and,
 $weight_decay = 0005$



FC Layer

Change the last fully connected layer into two.



Train one Epoch

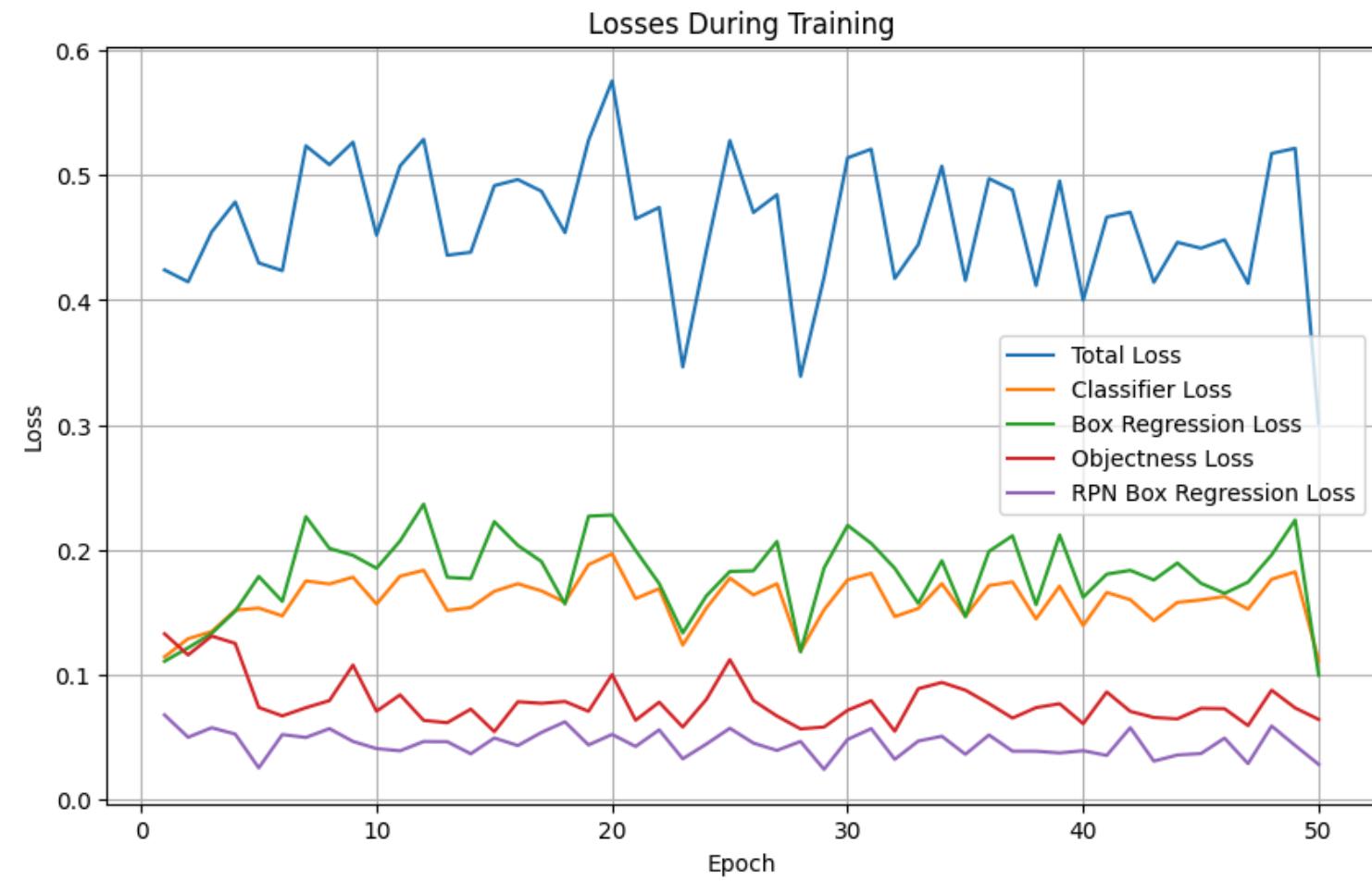
We use the train one epoch function from engine.py



Scheduler

Decrease the learning rate to optimize the model

The Training Performance (Loss)



The graph shows losses approximately above 0.05 or 5%, meaning the model has a high loss in creating the boxes for person detection and needs to be enhanced.

The loss scores consist of the following:

- Total loss: the total of all losses,
- Classifier loss: the loss that calculates the boxes over the wrong objects,
- Box regression loss: this loss measures how well the predicted bounding box matches the ground truth bounding box,
- Objectness loss: this loss is used to determine whether the proposed regions from the Region Proposal Network (RPN) actually contain an object (foreground) or not (background),
- RPN box regression loss: this loss measures how well the box proposals from RPN match the ground truth.

The Training Performance (AP and AR)

Average Precision (AP) @[IoU=0.50:0.95 area= all maxDets=100] = 0.158
Average Precision (AP) @[IoU=0.50 area= all maxDets=100] = 0.404
Average Precision (AP) @[IoU=0.75 area= all maxDets=100] = 0.082
Average Precision (AP) @[IoU=0.50:0.95 area= small maxDets=100] = 0.035
Average Precision (AP) @[IoU=0.50:0.95 area=medium maxDets=100] = 0.258
Average Precision (AP) @[IoU=0.50:0.95 area= large maxDets=100] = 0.303
Average Recall (AR) @[IoU=0.50:0.95 area= all maxDets= 1] = 0.118
Average Recall (AR) @[IoU=0.50:0.95 area= all maxDets= 10] = 0.254
Average Recall (AR) @[IoU=0.50:0.95 area= all maxDets=100] = 0.375
Average Recall (AR) @[IoU=0.50:0.95 area= small maxDets=100] = 0.204
Average Recall (AR) @[IoU=0.50:0.95 area=medium maxDets=100] = 0.405
Average Recall (AR) @[IoU=0.50:0.95 area= large maxDets=100] = 0.505

Similar to the previous result, we have a small score for precision and recall, meaning the model needs to be enhanced. The table shows that the highest precision is at IoU = 0.5, indicating that the model mainly yields the box at 50% of the actual box. The recall is a bit higher than the precision for IoU = 0.50:0.95, which tells us that the model is better at avoiding predicting a person as a non-person than vice versa. If we look at the image size, the model is worse at predicting small objects compared to medium and large objects.

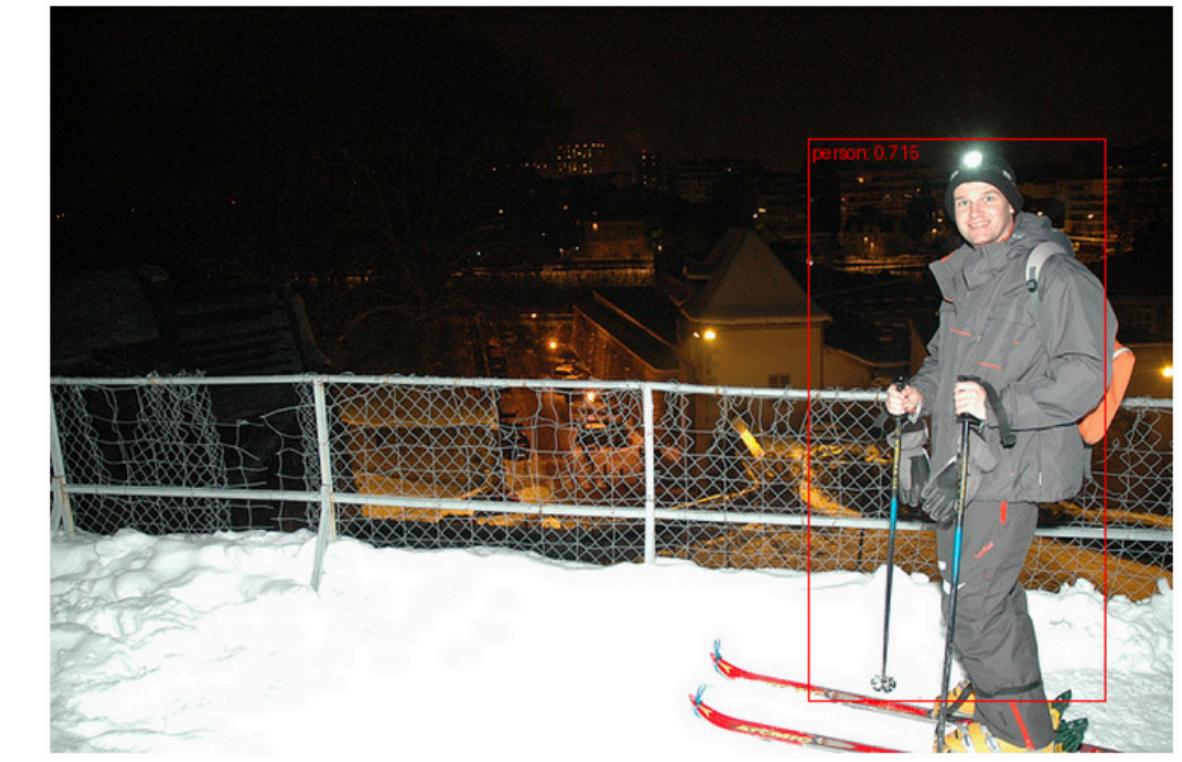
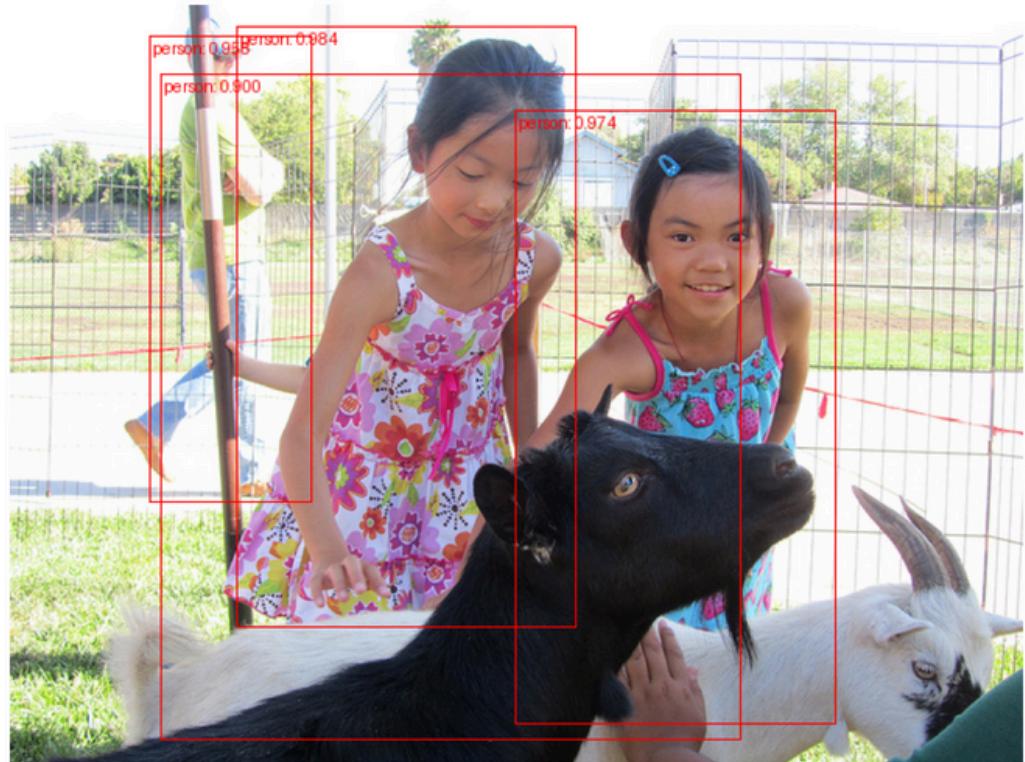
Model Testing

- Likewise, over the test dataset, the highest MAP score is when the threshold is 0.5, meaning that the model is better at predicting the bounding box at IoU greater or equal to 50%.
- However, the higher the IoU threshold, the worse the MAP score will be, and thus, the prediction will get worse.
- This indicates that the model can predict the object's position but cannot create a bounding box that covers the whole object at a high percentage.
- The inference latency is small, showing that the model runs quite fast.

MAP Score	
Threshold	Score
0.5	0.4632
0.6	0.3746
0.7	0.1675
0.8	0.0910
0.9	0.0395

Inference Latency = 0.1163 seconds

Model Visualization



From the first picture, we can see that there is a bounding box so large that it covers more than one object. The model sometimes also predicts a non-person as a person. A shadow in the second image, for instance, is detected as a person, while it is not. The last picture is better, but still, the box does not cover the whole person but instead covers outside of the person.

YOLO

Model Preparation and Construction

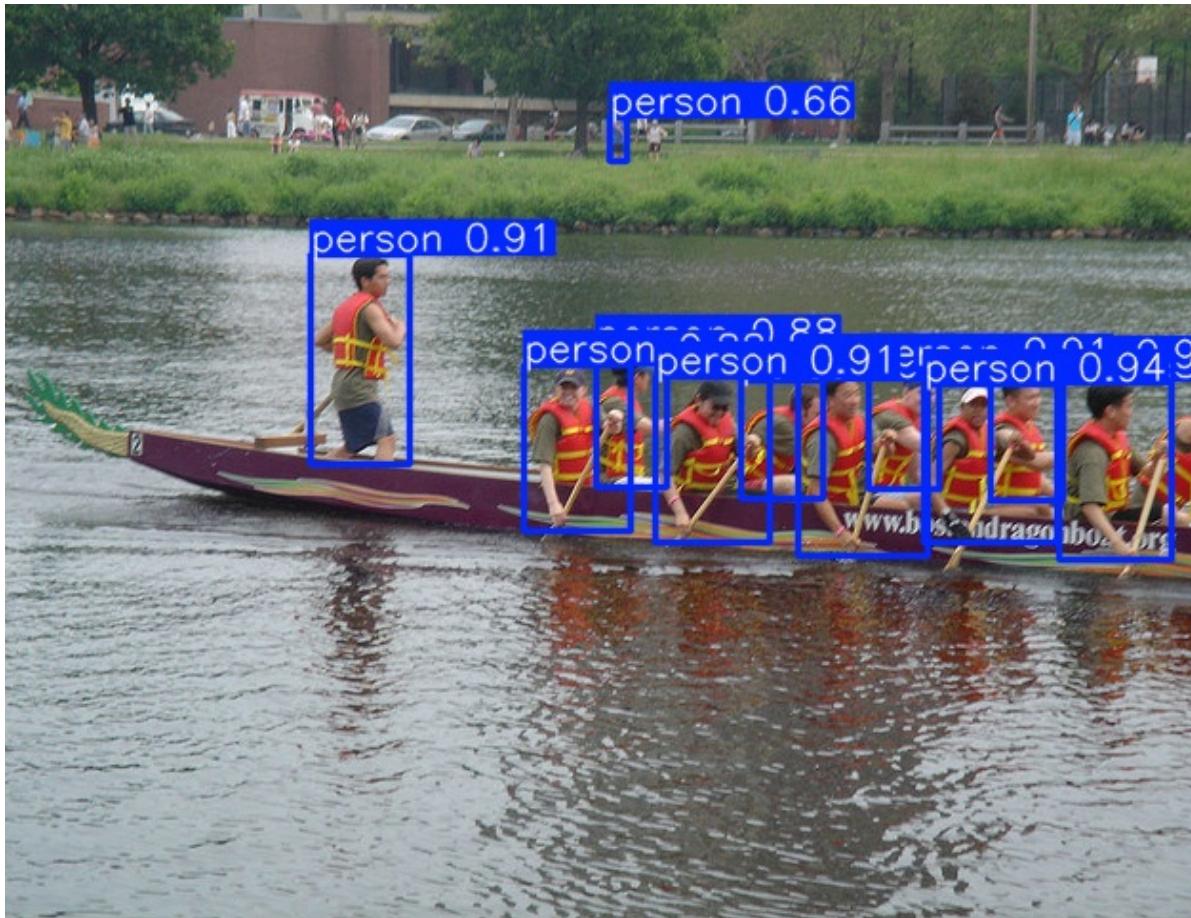
- After loading the images, we changed the box format to YOLO format, that is (X centre, Y centre, Width, and Height).
- Then, we created the YAML file, which consists of the root path, train path, val path, test path, number of classes, and class name.
- We import the YOLO function from Ultralytics and set the function to task=detect mode=train model=yolov8x.pt data=coco.yaml epochs=100.
- The function will automatically save the best resulting model to the runs folder.

Model Result

We obtain the following result for this model.

Metric	Score
Box Loss	0.6045
Classification Loss	0.4289
MAP 50	0.588
MAP 50:95	0.354

Model Visualization



This image showcases a YOLO (You Only Look Once) model output for object detection. The model has identified multiple people in a dragon boat race, with each person surrounded by a bounding box labeled as "person" along with a confidence score (ranging from 0.66 to 0.94). These confidence scores indicate how certain the model is about each detection.

Comparison

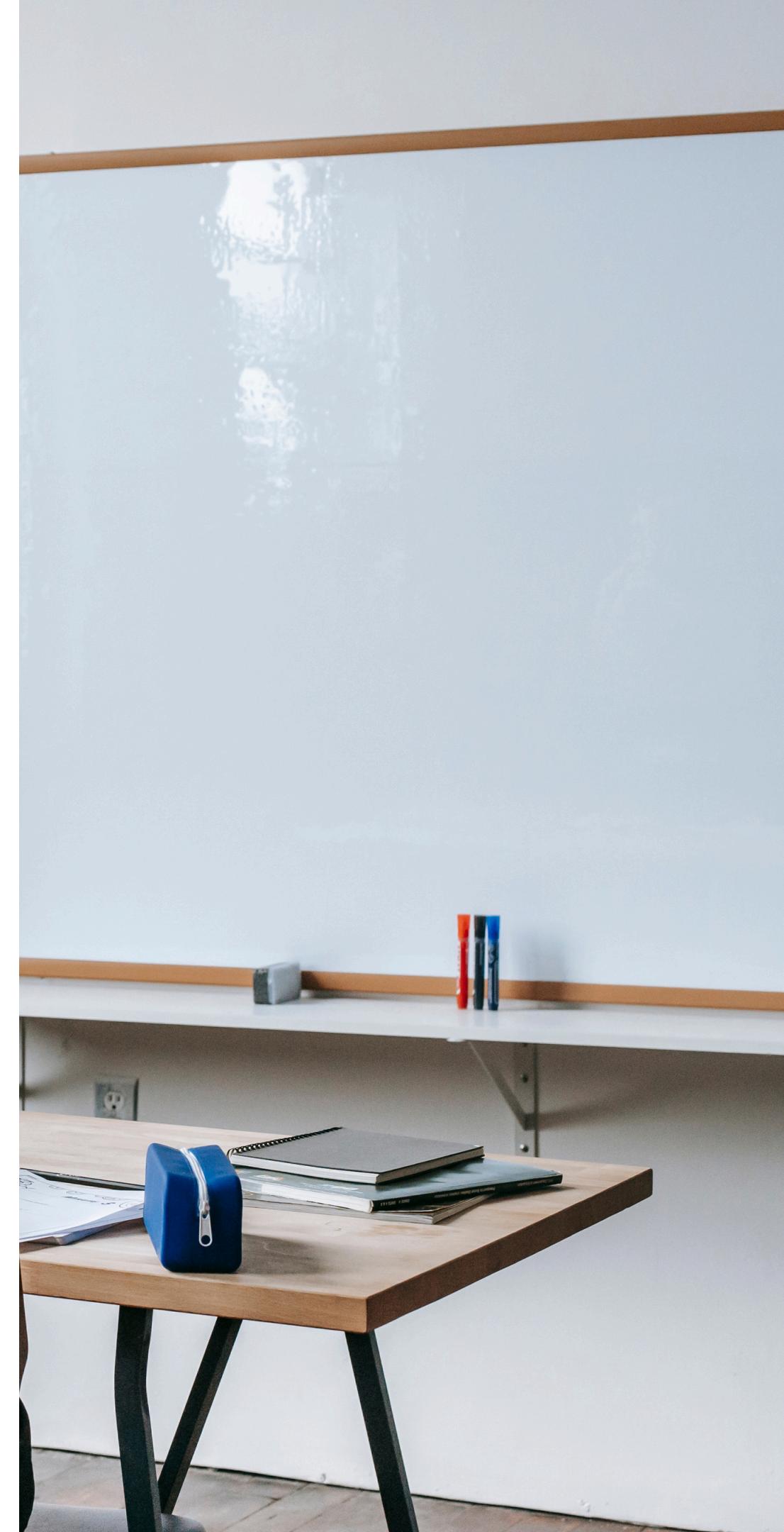
Metric	Faster RCNN	YOLO
Box Loss	0.0788	0.6045
Classification Loss	0.0994	0.4289
MAP 50	0.4632	0.588

The results show that although Faster RCNN obtained a small loss value during the training process, its MAP is still lower than YOLO. Hence, we decided that YOLO is better in this study case.

Future Improvement

The model needs improvements to obtain the best performance based on the losses and MAP scores. The enhancement can be:

- Since in this project, we only use 500 images, to escalate the performance, we can increase the number of images,
- The epoch can also be inclined as we only use 50 epochs,
- The learning rate can be considered to be altered,
- We can try to use other optimizers than SGD such as Adam optimizer as well as its parameters.



Real-world Applications

The person tracking can be applied to real-world situations such as:

1. Tracking suspicious people visiting our house or office to prevent robbery or any dangerous situations,
2. Detecting crowds in an event or public spaces to ensure safety and security,
3. Assisting autonomous vehicles to detect pedestrians and take action to avoid accidents, as well as to improve traffic light timings and road safety,
4. Tracking players in a sports competition so that it can track athletes during games for performance analysis, training, and enhancing viewer experiences,
5. It can be combined with a drone to detect people in a dangerous situation, such as during a natural disaster, so that they can be evacuated as soon as possible,
6. We can implement it for fall detection, especially for babies and hospital patients, to avoid accidents.

Conclusion

Person tracking is one of the object detection problems, and it has many implementations in the real world, such as CCTV for security and traffic and Autonomous Vehicle for pedestrian detection. There are some models that can be applied for person tracking, and two of them are Faster RCNN and YOLO.

In this study case, we use the COCO dataset to train both models for person tracking. Based on the previous performance and result, we conclude that the best model for this case is YOLO as it yields higher MAP.





Thank You !