

# Norfolk CV DL Homework

Vibhav Nirmal

Date: 06/03/2023

## Contents

1. Introduction.....	1
2. Model Training.....	2
<b>Data Preparation .....</b>	<b>2</b>
<b>Model configuration and Training.....</b>	<b>3</b>
<b>Model Evaluation .....</b>	<b>4</b>
3. Future Improvement.....	5
Part 1: How to run code .....	6
Part 2: Knuckle Pin detection .....	7
References:.....	9

## 1. Introduction

In this report, I present my approach to modeling a coupler detector using the provided data and the guidelines in the pdf document.

I considered two choices of object detection algorithms after looking at the dataset:

- **Faster RCNN**, which is a two-stage algorithm that first proposes potential regions of interest using a region proposal network (RPN) and then classifies and refines these proposals, and
- **YOLO**, which is a one-stage object detection algorithm that provides real-time detection capabilities and is faster than RCNN for edge device computation.

Considering that the primary objective was object detection with edge devices, I opted to utilize the YOLO algorithm.

I implemented the YOLO model in Python and PyTorch. I trained it on the data that was already in YOLO format, which made the task easier. The model takes an image as input and outputs a bounding box with a class label for the detected coupler.

## 2. Model Training

### Data Preparation

The dataset consisted of 200 images with bounding box annotations for the couplers. The images had different widths but the same height of 1024 pixels. I divided the dataset into **180 training images and 20 validation images** for the model training. Several **image processing parameters** were used for data augmentation. Some of them are listed below.

- **Image rotation:** 0.0 degrees (+/- deg) since all the images depicted the train and camera aligned horizontally.
- **Image translation:** A translation value of **0.2 (+/- fraction)** was applied to introduce slight variations in the object's position within the image.
- **Image scale:** **0.5 (+/- gain)**, allowing for both zoom-in and zoom-out effects on the coupler region.
- **Image perspective:** Since the provided dataset showed consistent perspective across the images, the perspective parameter was set to 0.0 (+/-fraction), maintaining the original perspective without any distortion.
- **Image Flip Left-Right:** To further diversify the training data, a **probability of 0.5** was assigned for flipping the images horizontally. This augmentation helps the model learn robustly from different orientations of the couplers.

All these image processing parameters are applied during data preparation by yolov7 model.



## Model configuration and Training

The YOLO model is known for its fast and efficient object detection performance, making it suitable for edge device computation. It is an **anchor-based model** that predicts potential bounding boxes and then applies **non maximal suppression** to get the best possible fit. To **optimize the model's performance**, I set the **initial learning rate** to  $lr0 = 0.01$  and **momentum** of 0.937 to **accelerate the learning process and facilitate faster convergence**. I also applied **Weight decay** of 0.0005 to **prevent overfitting and promote generalization**.

- Model was initialized with pre-trained weights provided by yolov7 (<https://github.com/WongKinYiu/yolov7/releases/download/v0.1/yolov7.pt>)
- Loss function:** YOLOv7 uses a composite loss function that consists of three components: a bounding box regression loss, a classification loss, and an objectness loss.
  - Lbox uses MSE for coordinates and dimensions, with Complete Intersection over Union as the overlap metric.
  - Lcls uses binary cross-entropy for multi-label classification.
  - Lobj uses binary cross-entropy with focal loss for bounding box confidence.
- Optimizer:** YOLOv7 uses an **SGD (Stochastic Gradient Descent) optimizer** with a cosine annealing learning rate schedule. This helps the model avoid getting stuck in local minima and achieve better convergence.
- I trained the model for **100 epochs**, with an **input dimension of 640x640** pixels and a total **batch size of 8**.

```

detect.py --weights yolov7.pt --source ./data --img-size 640 --conf-thres 0.25 --iou-thres 0.45 --device cuda --save-txt --save-conf --nosave --project runs/detect --name defectDataCustom --exist-ok

```

Epoch	gpu_mem	box	obj	cls	total	labels	img_size
38/99	5.94G	0.01347	0.001159	0.001769	0.01694	12	640: 100%
		Class	Images	Labels	P	R	mAP <sub>0.5</sub> mAP <sub>0.5:0.95</sub> mAP <sub>0.5:0.95</sub>
		all	18	19	0.554	0.894	0.845 0.392
39/99	5.94G	0.03048	0.00710	0.001652	0.04084	9	640: 100%
		Class	Images	Labels	P	R	mAP <sub>0.5</sub> mAP <sub>0.5:0.95</sub> mAP <sub>0.5:0.95</sub>
		all	18	19	0.66	0.889	0.924 0.429
40/99	5.94G	0.04245	0.00726	0.001825	0.053	14	640: 100%
		Class	Images	Labels	P	R	mAP <sub>0.5</sub> mAP <sub>0.5:0.95</sub> mAP <sub>0.5:0.95</sub>
		all	18	19	0.833	0.934	0.981 0.537
41/99	5.94G	0.03254	0.00835	0.001567	0.04245	13	640: 100%
		Class	Images	Labels	P	R	mAP <sub>0.5</sub> mAP <sub>0.5:0.95</sub> mAP <sub>0.5:0.95</sub>
		all	18	19	0.944	0.878	0.97 0.442
42/99	5.94G	0.04123	0.00803	0.001907	0.05204	21	640: 100%
		Class	Images	Labels	P	R	mAP <sub>0.5</sub> mAP <sub>0.5:0.95</sub> mAP <sub>0.5:0.95</sub>
		all	18	19	0.739	0.889	0.967 0.545
43/99	5.94G	0.04189	0.00825	0.001793	0.0521	15	640: 100%
		Class	Images	Labels	P	R	mAP <sub>0.5</sub> mAP <sub>0.5:0.95</sub> mAP <sub>0.5:0.95</sub>
		all	18	19	0.797	0.943	0.959 0.548

## Model Evaluation

The trained model was evaluated on a separate validation dataset with 20 images to measure its performance metrics, such as precision, recall, and mean average precision (mAP)

**With 180 train - 20 val split**, I got the following results: (For more information I have shared the results.txt)



- The model achieved a **mean average precision (mAP) of 0.9962**, which means that it **correctly detected 99.62%** of the objects in the images.
- The model also achieved **high precision (1), recall (1) and F1-score**, which means that it **had very low (zero) false negatives and false positives**.
- The model also had **very low losses for bounding box regression (0.009411), objectness prediction (0.002513), and class prediction (0)**

I did run the model on all 459 images and almost all the images had been correctly classified visually. As those images were not labeled, I labeled 20 images from that dataset (using `labelImg`) and evaluated the model.

Testing on **20 newly annotated test images** produced **mAP of 0.995**.

```
(venv) PS E:\coding\Norfolk\yoloV7-main> python test.py --weights runs\tain\norfolk-yolo2\weights\best.pt --data data\data\trainDataCustom.yaml --task test --name yolo_det_test_couppler  
Namespace(weights=['runs\\tain\\norfolk-yolo2\\weights\\best.pt'], data='data\\trainDataCustom.yaml', batch_size=32, img_size=640, conf_thres=0.001, iou_thres=0.65, task='test', device='', single_cls=False, augment=False, verbose=False, save_txt=False, save_hybrid=False, save_conf=False, save_json=False, project='runs/test', name='yolo_det_test_coupler', exist_ok=False, no_trace=False, v5_metric=False)  
YOLOR 2023-6-2 torch 1.12.1+cu113 CUDA:0 (NVIDIA GeForce RTX 3070 Ti Laptop GPU, 8191.5MB)  
  
Fusing layers...  
RepConv.fuse_repvgg_block  
RepConv.fuse_repvgg_block  
RepConv.fuse_repvgg_block  
IDetect.fuse  
E:\\coding\\Norfolk\\yolv7-main\\.venv\\lib\\site-packages\\torch\\functional.py:478: UserWarning: torch.meshgrid: in an upcoming release, it will be required to pass the indexing argument. (Triggered internally at ..\\aten\\src\\ATen\\native\\TensorShape.cpp:2895.)  
    return _VF.meshgrid(tensors, **kwargs) # type: ignore[attr-defined]  
Model Summary: 314 layers, 36481772 parameters, 6194944 gradients, 103.2 GFLOPS  
Convert model to Traced-model...  
Traced script module saved!  
model is traced!
```

	Class	Images	Labels	P	R	mAP@.5	mAP@.5:.95:
Test Results	all	20	20	1	1	0.995	0.512

Speed: 22.8/7.8/30.6 ms inference/NMS/total per 640x640 image at batch-size 32  
Results saved to runs\test\yolo\_det\_test\_coulper3  
( .venv ) PS E : \ coding \ Norfolk \ yolov7 - main >

### 3. Future Improvement

If I had infinite time and/or resources, I would try to improve the model performance by using:

- **more data:** More images from **different sources, conditions, and perspectives** to **increase the generalization** ability of the model (at least 1000 images per class)
- **different architectures:** Experiment with different **backbone, neck, and head architectures** to find the **optimal trade-off between speed and accuracy**.
- **different models:** RCNN, UNet or any other object detection/semantic segmentation models
- **different hyperparameters** and **different evaluation metrics**

## Part 1: How to run code

I used **PyTorch** as the programming platform and trained the model on a **3070Ti GPU**. The code (**detectTrainCoupler.py**) and the model file (**couplerbest.pt**) are attached in this folder.

The model requires Python version 3.8 or higher to run. Use the following command to install dependencies. (Recommended to use virtualenv)

```
pip install -r requirements.txt
```

(Note that I have used torch 1.12.1 with cuda 11.3. Make sure to change according to your requirements)

To run the inference on an image and save in a file,

```
python detectTrainCoupler.py --img-path <path_to_image> --weights  
<path_to_model>
```

Following command can be used to run inference on multiple images in a folder:

```
python detectTrainCoupler.py --folder-path <path_to_folder_containing_images>  
--weights <path_to_model>
```

This will process all the images in the folder and save the output as images or a video file (10 fps)

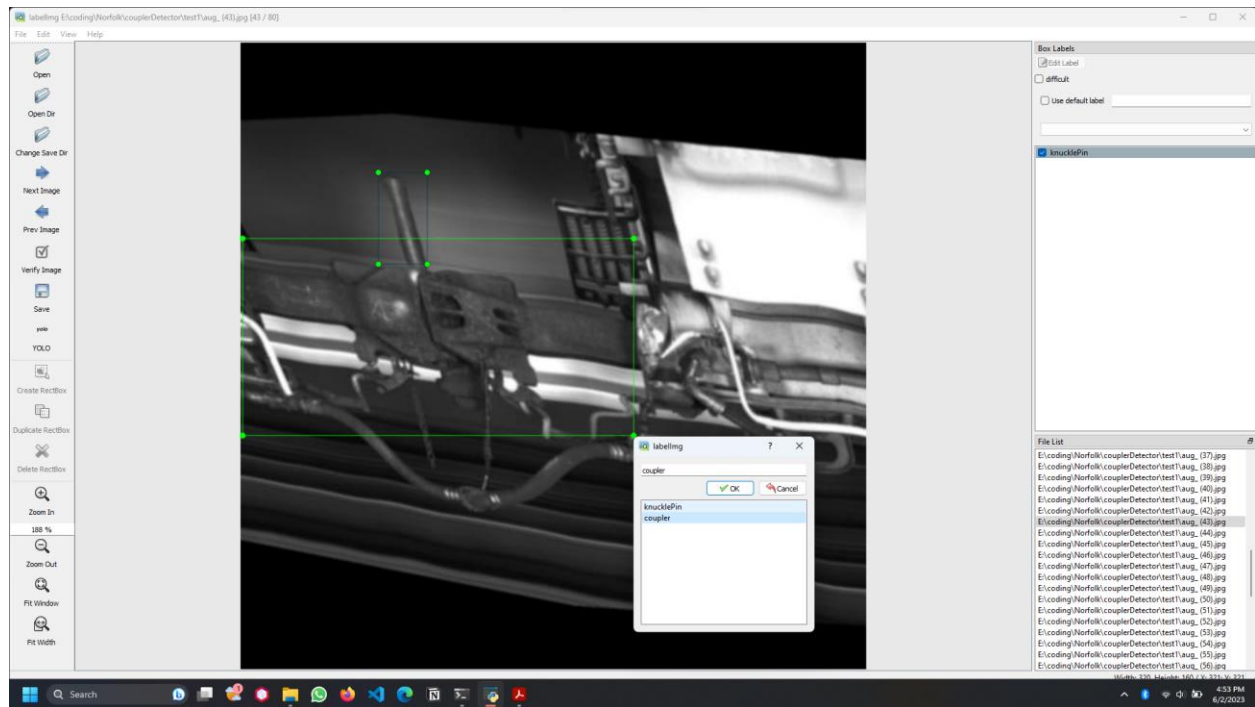
### Example:

```
python detectTrainCoupler.py --weights couplerbest.pt --folder-path testImgs  
--save-as-video False --conf-thres 0.4 --verbose False
```

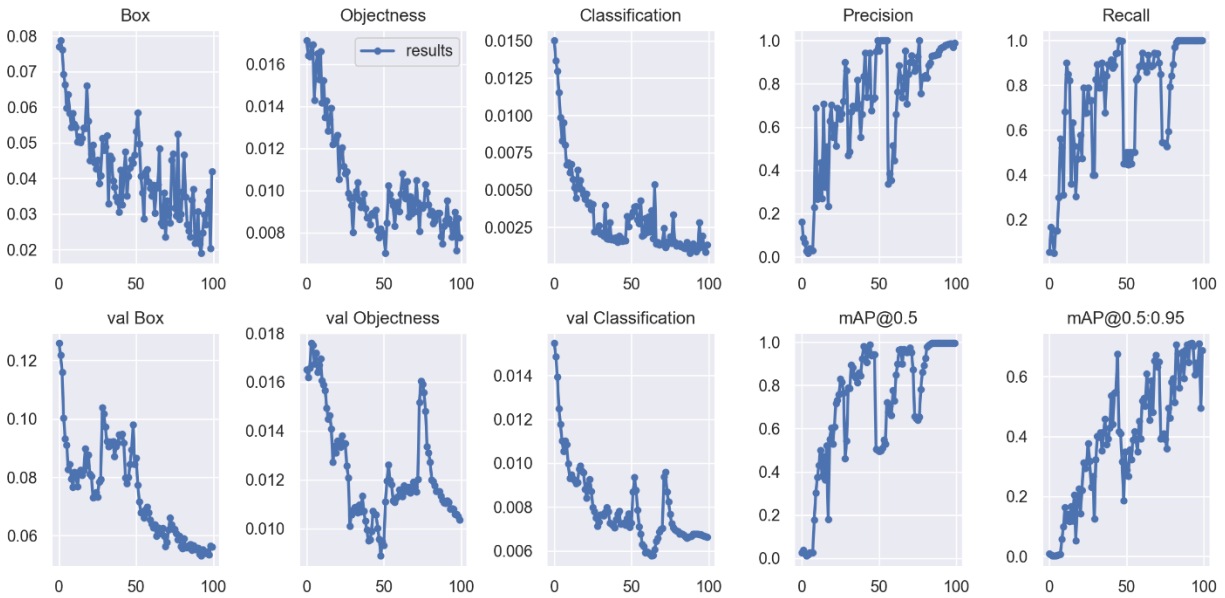


## Part 2: Knuckle Pin detection

Since I only had one image to detect defects, I used some data augmentation techniques to improve the model's performance. I have augmented it by various means to create **60 training images** and **10 validation images**.



However, my model was not able to predict well. I have attached the results and trained model.



To detect the defect by having a single image at hand, I can fine-tune my localizer model (couplerbest.pt) that can detect the coupler and Knuckle Pin using transfer learning and data augmentation. For yolo, this involves:

- Creating a dataset with labeled images of the defect and applying transformations such as rotation, scaling, cropping, etc. to generate more images.
- Modifying the config file to match the new classes and anchors.
- Using couplerbest.pt as the initial weights
- Training and tuning the model on the new dataset.
- Evaluating the model on a validation or test set

**Data augmentation** and **transfer learning** are techniques that can help overcome data scarcity by increasing the diversity of the training dataset and leveraging the knowledge of a pre-trained model. I can use yolov7 as my pre-trained model and fine-tune it on my new dataset that contains augmented images of the defect.

Also, I can use **simulated images** to train the model to make it robust and generalized. Having different weather in simulation like rain, snow or mud can help model make better predictions over unseen images.



## References:

YOLOv7 is the primary source of the whole project.

[YOLOv7: The Fastest Object Detection Algorithm \(2023\) - viso.ai](#)

[YOLOv7: A deep dive into the current state-of-the-art for object detection | by Chris Hughes & Bernat Puig Camps | Towards Data Science](#)