
Recent Advances in Machine Learning - Object Detection and Robot Control

University of Siegen

Shantanu Shirsath Aakash Deshpande
1699551 1701953

Abstract

Object detection is one of the most prominent tasks in the field of computer vision that involves identifying and localizing objects present in a frame of image or video. This is a crucial method which finds applications in fields like autonomous driving, video surveillance, medical imaging, etc.

In this report, we focus on the implementation of real-time object detection and control of a robot. The project features integration of deep learning knowledge and hardware setup along with control engineering. Each segment of this project is well explained with the help of graphs or images.

1 Introduction

The main objective of this project is to enable a Fischertechnik robot to detect and chase a specific object, a Cat in this case. A pre-trained deep learning model is retrained using a custom made dataset and deployed on a Raspberry Pi. The robot uses a webcam and the model to detect and track a cat in real-time. PID controllers compute motion commands based on the cat's position and size, enabling responsive and continuous pursuit.

We have combined various datasets to form a single integrated dataset that is well balanced, keeping in mind the deep learning aspect or retraining of a model. The topic of Dataset generation will put more light on this.

Selecting a deep learning model is a crucial task as the efficiency of the whole system depends on how well the model is able to predict or identify the object. Reasons for selection of YOLOv8n[1] and retraining will be discussed in the subsequent sections.

Once the model is trained and a checkpoint point is generated, it needs to be integrated on a hardware, in this case, a Raspberry Pi-5. The conversion methods and the lightweight inference engine selected are very important to make the model more robust. We have selected a high performance neural network inference framework named NCNN developed by Tencent. The details of which are discussed further.

Once the hardware is able to the real time inference and is able to classify the objects captured by the robot's webcam, the control aspect of this project comes into play. We have used 2 simple PID controllers to control the motion of the robot based on the errors generated. More about the Control engineering is mentioned in the section Control Engineering.

2 Dataset generation

For any Deep learning project which deals with training or retraining, the data plays a major role. A right dataset can help save countless hours spent fine tuning the model to increase its efficiency (collection of numerous parameters). Since we want our robot to detect a Cat, we have generated a mega-dataset of Cats. In initial stage, the dataset was not balanced and many issues were observed. The following table 1, lists the problems and action taken to solve that problem and table 2 shows how the dataset was distributed and balanced in the old and new dataset.

Problem	Solution
Model assumed a cat is always present.	Added non-cat images to reduce overfitting.
High false positives on cat-like objects.	Balanced cat and non-cat ratio to 50-50.
Model detected only one cat breed.	Added diverse breeds and custom cat images.

Table 1: Summary of key issues in the old dataset and corrective actions taken.

Source	Category	Old Dataset		New Dataset	
		Share (%)	Count	Share (%)	Count
COCO[2]	Non-cat	10%	478	50%	8,794
COCO	Cat	90%	4,298	25%	4,298
Roboflow[3]	Cat	–	–	14%	2,433
Custom	Cat	–	–	11%	2,063
Total		100%	4,776	100%	17,588

Table 2: Comparison of old and new dataset distributions used for cat vs. non-cat image classification.

37 The generation of a custom dataset included recording a video and extracting the frames. Then each
38 frame was augmented as mentioned below.

- 39 • **A.OneOf([...], p=1.0)**
40 Applies exactly one of the following three transformations (with 100% probability):
- 41 – **A.HorizontalFlip(p=1.0)**: Flip the image left to right.
 - 42 – **A.VerticalFlip(p=1.0)**: Flip the image top to bottom.
 - 43 – **A.Rotate(limit=30, p=1.0)**: Rotate the image at random angle between -30° & $+30^\circ$.
- 44 • **A.RandomBrightnessContrast(p=0.5)**
45 With 50% probability, randomly changes:
- 46 – **Brightness**: Makes the image lighter or darker.
 - 47 – **Contrast**: Increases or decreases the difference between light and dark areas.
- 48 • **A.HueSaturationValue(p=0.5)**
49 With 50% probability, randomly adjusts:
- 50 – **Hue**: Adjusts the color tone.
 - 51 – **Saturation**: Modifies the intensity of colors.
 - 52 – **Value**: Alters the brightness component in HSV color space.

53 After the augmentation, each image was annotated using inference of YOLOv8n itself and
54 incorporated in the main dataset.

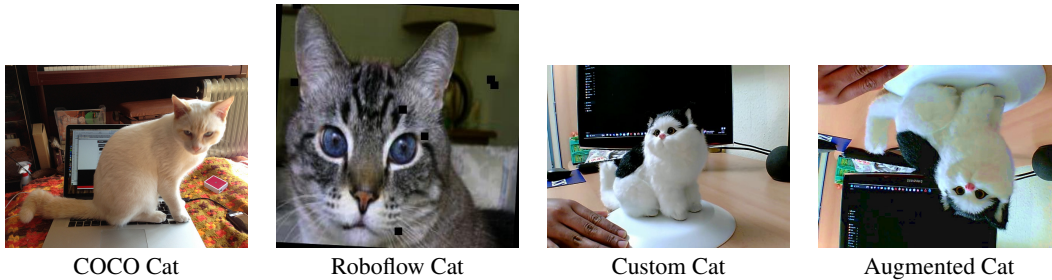


Figure 1: Sample cat images from various datasets

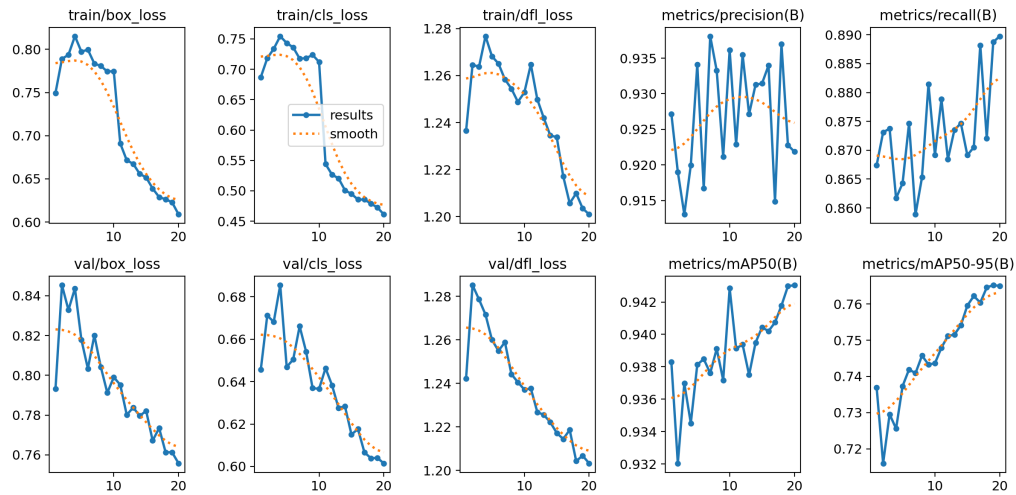
3 Selecting, training and Implementation of the model onto a Raspberry Pi-5

3.1 Selection and training

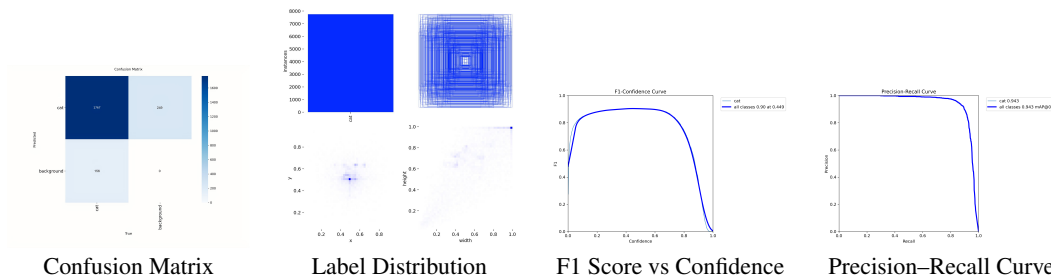
The selection of the model was dependent on the following things

- The architecture, model size and training simplicity.
- The Accuracy and speed of the inference.
- The Simplicity of exporting it onto Raspberry Pi.

After comparing the models like YOLOv5, YOLOv8n, SSD(with mobilenet), Faster R-CNN, EfficientDet and RetinaNet, we concluded that YOLOv8n is the best model to use for our application. Since the Model is already taught to extract the features on a gigantic dataset, there was no need to retrain it completely with our small dataset. The training was carried out by freezing the Backbone and Neck of the model which form the first 21 layers. The head which is the 22nd layer is responsible for classifying the objects and we fine turned the model to classify only 1 class i.e Cat. After splitting the data into train and validation (4:1), the training was carried out in 2 stages. First for 30 epochs and then for 20 more epochs. The results generated by the YOLO are as shown below:



Training Curves



Confusion Matrix

Label Distribution

F1 Score vs Confidence

Precision-Recall Curve

3.2 Implementation on Raspberry Pi 5

For evaluating inference performance on Raspberry Pi, we tested four popular model formats with inference using the ultralytics library: **PyTorch (.pt)**, **TorchScript**, **ONNX**, and **NCNN**. Below is a brief overview of each:

- **PyTorch (.pt)**: Native PyTorch model format used during training and initial evaluation.
- **TorchScript**: A serialized and optimized representation of a PyTorch model, enabling deployment in C++ environments.

- **ONNX (Open Neural Network Exchange):** An open format designed to allow interoperability between different deep learning frameworks.
- **NCNN:** A high-performance neural network inference framework optimized specifically for mobile and edge devices, especially those using ARM processors.

Inference Performance on Raspberry Pi

We benchmarked each format on four key metrics: preprocessing time, inference time, post-processing time, and achieved FPS (frames per second). The results are summarized in Table 3.

Table 3: Inference Performance Comparison

Format	Preprocess (ms)	Inference (ms)	Postprocess (ms)	FPS
TorchScript	5–8	456	0.5–1.5	1.5–2
PyTorch (.pt)	5–6	350	0.5–1	2.5–3
ONNX	5	300	0.5–1	3–3.5
NCNN	5–6	110	1–2	8–9

From the above comparison, it is evident that the NCNN format significantly outperforms others in terms of inference speed and achieved FPS on Raspberry Pi. This is primarily due to NCNN’s design, which is optimized for ARM-based architectures and lightweight inference on edge devices. Therefore, for real-time deployment on Raspberry Pi, **NCNN was chosen as the preferred format** due to its efficiency and low-latency performance.

4 Control Engineering

We control two axis motion of the Robot x-Axis (lateral motion) and y-axis (longitudinal motion). In order for the robot to chase the cat, the control problem can be defined so as to minimize the distance between robot and cat along both x and y axis. The motion along x-axis defines the lateral distance or the direction in which the robot points while the motion along the y axis defines the heading direction of the robot.

4.1. Motion Along the X-Axis

To control the motion along x-axis we take help of the bounding box center. We calculate the bounding box center and try to minimize the distance between bounding box center and the frame center of the camera feed. This helps the robot to adjust its direction towards the object trying to keep the object always at the center of the frame.

4.1.1 lateral controller

To minimize the lateral distance between the center of the camera frame and the center of the bounding box (i.e., the cat’s position), we first compute the error:

$$e(t) = x_{\text{center}}^{\text{frame}} - x_{\text{center}}^{\text{bbox}} \quad (1)$$

where:

- $x_{\text{center}}^{\text{frame}}$ is the x-coordinate of the center of the camera frame,
- $x_{\text{center}}^{\text{bbox}}$ is the x-coordinate of the center of the detected bounding box.

This error $e(t)$ represents how far off-center the object (cat) is in the horizontal direction. To correct this error and steer the robot to keep the object centered, we use a Proportional-Derivative (PD) controller.

A PD controller generates a control signal based on the present error and the rate of change of the error. The control law for a PD controller is given by:

$$u(t) = K_p \cdot e(t) + K_d \cdot \frac{de(t)}{dt} \quad (2)$$

110 where:

- 111 • $u(t)$ is the control signal (e.g., angular velocity or steering adjustment),
- 112 • K_p is the proportional gain,
- 113 • K_d is the derivative gain,
- 114 • $e(t)$ is the current error,
- 115 • $\frac{de(t)}{dt}$ is the derivative of the error with respect to time.

116 The proportional term ($K_p \cdot e(t)$) provides a correction proportional to the error, encouraging the
 117 robot to move such that the object is at the center of the frame. The derivative term ($K_d \cdot \frac{de(t)}{dt}$) helps
 118 dampen the response and reduce overshooting by considering how quickly the error is changing.

119 The computed control signal $u(t)$ is then sent to the motor controller to adjust the robot's orientation,
 120 ensuring that the object remains centered in the camera frame.

121 If we assume that the sampling time Δt is constant, we can approximate the derivative term using the
 122 finite difference method:

$$\frac{de(t)}{dt} \approx \frac{e(t) - e(t-1)}{\Delta t} \quad (3)$$

123 Substituting this into the PD control law:

$$u(t) = K_p \cdot e(t) + K_d \cdot \frac{e(t) - e(t-1)}{\Delta t} \quad (4)$$

124 Since Δt is constant, we can define a new derivative gain:

$$K'_d = \frac{K_d}{\Delta t} \quad (5)$$

125 Thus, the simplified discrete-time control law becomes:

$$u(t) = K_p \cdot e(t) + K'_d \cdot (e(t) - e(t-1)) \quad (6)$$

126 This form eliminates the explicit time derivative, making it more suitable for digital implementation
 127 on an embedded system.

128 4.1.2. Motion Along the Y-Axis

129 To control the motion along the y-axis (i.e., forward or backward movement), we rely on the area
 130 of the bounding box detected around the object (the cat). Assuming the cat has a roughly constant
 131 physical size, a smaller bounding box area indicates that the cat is farther from the robot, while a
 132 larger area implies proximity.

133 We define a target area A_{target} that represents 80% of the total camera frame area:

$$A_{\text{target}} = 0.8 \cdot A_{\text{frame}} \quad (7)$$

134 The control objective is to maximize the bounding box area A_{bbox} such that it approaches A_{target} . To
 135 do this, we define the error in area as:

$$e(t) = A_{\text{target}} - A_{\text{bbox}}(t) \quad (8)$$

136 A positive error means the object is still too far away. To minimize this error, we apply a **Proportional**
 137 **(P) controller**, which adjusts the control signal based on the current error:

$$u(t) = K_p \cdot e(t) \quad (9)$$

138 where:

- 139 • $u(t)$ is the control signal (linear velocity sent to the motors),
- 140 • K_p is the proportional gain,
- 141 • $e(t)$ is the current area error.

142 The control signal $u(t)$ determines the forward motion of the robot. As the robot approaches the
 143 object and the bounding box area increases, the error decreases, causing the robot to slow down and
 144 helping it maintain a desired distance from the object.

145 5 Controller tuning

146 5.1 Lateral PD Controller Behavior

147 The PD controller for lateral motion was tuned by varying the derivative gain K_d while keeping the
 148 proportional gain K_p fixed. As shown in Figure ??, increasing K_d initially improves response time
 149 by reducing overshoot. However, higher values also lead to increased oscillations and instability due
 150 to the derivative term's sensitivity to noise and sharp error changes.

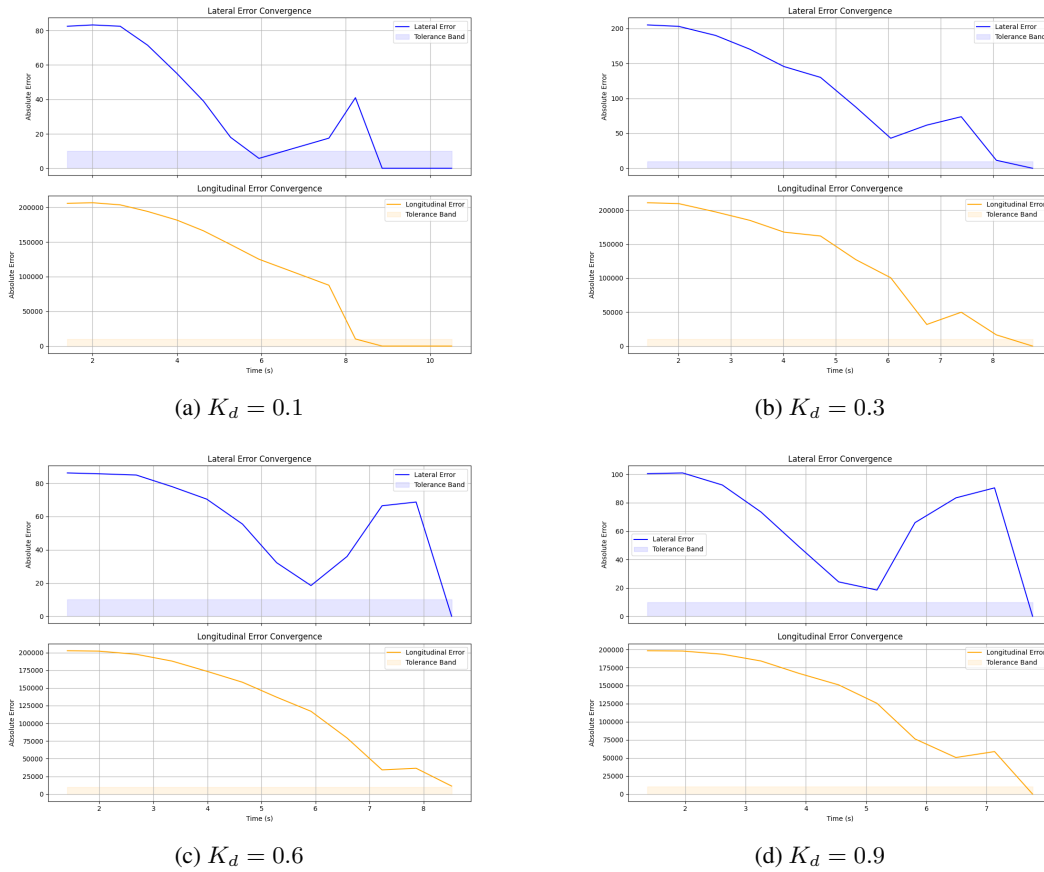


Figure 2: Comparison of Lateral and Longitudinal Error Convergence for varying K_d in the PD controller

151 As seen in the plots, a lower K_d results in smooth and stable convergence, whereas higher K_d values
152 introduce fluctuations and delay in settling within the tolerance band.

153 5.2 Longitudinal P Controller Behavior

154 The longitudinal controller, being purely proportional, exhibits a monotonically decreasing error
155 profile across all experiments. Its response is inherently slower but stable, as shown in all subplots of
156 Figure 2. The absence of derivative or integral action avoids oscillations but limits responsiveness.

157 5.3 Conclusion

158 From the tuning experiments:

- 159 • Low to moderate K_d values offer a good trade-off between responsiveness and stability in
160 PD control.
- 161 • High derivative gain can lead to undesirable oscillations and overshoots due to noise ampli-
162 fication.
- 163 • The longitudinal P controller provides stable convergence, but the system could benefit from
164 adding derivative or integral components for faster performance. But due to the limitation of
165 the responsiveness of our hardware we limit our controller to simple p controller.

166 References

- 167 [1] Glenn Jocher, Jing Qiu, and Ayush Chaurasia. Ultralytics yolo, 2023.
- 168 [2] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
169 Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet,
170 Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*,
171 pages 740–755, Cham, 2014. Springer International Publishing.
- 172 [3] Mohamed Traore. Cats dataset. [https://universe.roboflow.com/](https://universe.roboflow.com/mohamed-traore-2ekkp/cats-n9b87)
173 mohamed-traore-2ekkp/cats-n9b87, nov 2022. visited on 2025-07-10.