

CV-BASED BATCH-BILLING SYSTEM FOR SUPERMARKET PRODUCTS USING YOLO

**(This project report has been submitted in partial fulfilment of the requirements
for the degree of Bachelor of Science in Electrical and Electronic Engineering)**



SUBMITTED BY

EXAM ROLL NO: 10356

REG. NO: 2016-315-002

SESSION: 2016-17

EXAM ROLL NO: 10313

REG. NO: 2016-615-027

SESSION: 2016-17

**DEPT. OF ELECTRICAL AND ELECTRONIC ENGINEERING
UNIVERSITY OF DHAKA**

October 10, 2021

Abstract

We propose a computer vision-based billing system for the check-out of retail items on the cash counter. Supermarkets employ traditional barcode scanners to detect retail products. The disadvantage of this prevailing check-out system is that it cannot detect a batch of products simultaneously and is hence time-consuming. The system also relies on the speed of the cashier's workflow in handling the cash counter. In this regard, our proposed system aims to make batch-billing of retail items faster and easier, providing a quality shopping experience for the customers. We have implemented a GUI-based check-out application that works on a YOLO-based object detection architecture to bill the retail products. With YOLO architecture at its core, the system scans retail items using a webcam placed above for detection. For training our model, we have chosen 16 local retail products of various categories and used both authentic and synthetic images of them as our dataset. Our system has achieved a reasonable real-time retail product detection accuracy to be implemented on the billing system paving the way for a fully automated supermarket experience in the future.

CONTENTS

Abstract	i
List of Figures	v
List of Tables	viii
1 INTRODUCTION	1
1.1 Motivation	1
1.2 Objectives	2
2 RELATED WORKS	3
2.1 ARC Vision-based auto retail	3
2.2 CV based Retail Store Product Detection	4
2.2.1 Block based Methods	5
2.2.2 Geometric Transformation based Methods	7
2.2.3 Saliency based Methods	7
2.2.4 Detector based Methods	8
2.2.5 User-in-the-loop Methods	9
2.3 A Computer Vision System Supporting Blind People	9
2.4 Toward Real-Time Grocery Detection for the Visually Impaired . .	12
2.5 Fine-Grained Grocery Product Recognition by One-Shot Learning .	14
3 THEORETICAL OVERVIEW	17
3.1 Hardware Specification	17
3.2 Deep learning and Neural Networks	18
3.3 Transfer Learning	20
3.4 Defining Object Detection	21
3.5 Object Detection Algorithms	24
3.5.1 R-CNN	24
3.5.2 SPP Net	25
3.5.3 Fast RCNN	25
3.6 YOLO Overview	27
3.6.1 History	27
3.6.2 How it works	28

3.6.3	YOLO Overview	30
3.6.4	YOLOv5	34
3.7	Data Augmentation	36
3.7.1	Basic Augmentation Techniques	36
3.7.2	Synthetic Data	37
3.8	YOLO Annotation	39
3.9	Optimizers	41
3.9.1	SGD (Stochastic gradient descent)	41
3.9.2	ADAM	43
3.10	YOLO Hyperparameters	44
3.11	Hyperparameter Evolution	45
4	METHODOLOGY	47
4.1	Dataset Description	47
4.1.1	Data Acquisition	50
4.1.1.1	Daraz Review Image Scrapping	50
4.1.1.2	Manual Photography	51
4.1.2	Data Cleaning	51
4.1.3	Pre-Processing	51
4.1.4	Data Annotation	52
4.1.5	Synthetic Data	53
4.1.5.1	Annotating Synthetic Data	55
4.1.6	Dataset Distribution	56
4.2	Model Overview	59
4.2.1	Hyperparameter Evolution	61
4.2.2	Training Description	62
4.3	Hardware Setup	63
4.4	GUI Implementation	64
4.4.1	Objectives of the Application	64
4.4.2	Description of the UI	65
4.4.3	Highlight Features of the App	66
4.4.3.1	Automatic Refresh	66
4.4.3.2	Batch Locking Mechanism	67
4.4.3.3	Receipt Generation	68
5	RESULT AND ANALYSIS	69
5.1	Dataset Analysis	69
5.2	Evaluation Metrics	71
5.2.1	Intersection Over Union (IoU)	71
5.2.2	Evaluation of Object Classification	72
5.3	Test setups	74
5.4	Models for Evaluation	77
5.5	Training Results	78
5.6	Results on Densely Arranged Objects	79

5.7	Results in Different Lighting Conditions	81
5.8	Analysis of Real-time detection	82
5.9	Analysis on Different Backgrounds	85
5.10	Analysis on Different Backgrounds	87
5.11	Effects of Hyperparameter Evolution	88
5.11.1	Values Obtained through evolution	89
5.11.2	Performance comparison on test sets	91
6	CONCLUSION AND FUTURE SCOPE	94
6.1	Discussion	94
6.2	Future Scopes	95
Bibliography		97
Appendix A: List of Acronyms		110

LIST OF FIGURES

2.1	Structure of the hardware used in the paper	4
2.2	Block representation of sliding window	6
2.3	Block representation of grid based method	6
2.4	Block representation of geometric transformation based method	7
2.5	Block representation of saliency-based method	8
2.6	Block representation of detector-based method	9
2.7	Sample images for different levels of abstract	10
2.8	Representation of Confusion matrices for the quaternary classification task using VGG, Resnet, Inception	11
2.9	Under-construction state of the mosaic represented by the green box in the FOV(Field of Vision) of the camera	13
2.10	Training (top) and testing (bottom) examples of (a) easy retail items, which gives zero false positives for a strict threshold; (b) hard retail items, which fails to detect with reasonable false positives	13
2.11	Proposed framework of the study	15
2.12	Production of candidate region by extracting feature map of the retail items on the shelf and comparing it with the feature points of the training sample	15
2.13	Generation of Attention Map Using SIFT	16
3.1	Model of Perceptron	19
3.2	Standard Deep Neural Network	19
3.3	Improvements due to transfer learning [1]	21
3.4	Roadmap of Object Detection [2]	22
3.5	Core Tasks of object detection	22
3.6	Comparison of semantic segmentation, classification and localization, object detection and instance segmentation. [3]	23
3.7	R-CNN Structure RCNN	24
3.8	Faster RCNN architecture	26
3.9	YOLO model with 7*7 grid was applied to input Image(Redmon, et al. 2016) [4]	28
3.10	Example of how to calculate box coordinates in a 448*448 image with S=3. [5]	29
3.11	Each grid cell makes B bounding box predictions and C class predictions. [5]	29
3.12	Matrix calculation for YOLO [6]	30

3.13	Network architecture of YOLO [7]	30
3.14	YOLOv1 model architecture. [5]	31
3.15	Building block of residual learning.	33
3.16	Darknet-53.	34
3.17	Image mixing via random cropping.[8]	37
3.18	Image composition via masking	38
3.19	Basic GAN Architecture [9]	38
3.20	Synthetic image generation with GAN (Yellow marked column is the closest training sample to the neighbouring output.)	39
3.21	Bounding Box	40
3.22	Local Minima for SVG algorithm	42
3.23	Saddle Point in minimizing the loss	43
3.24	Hyperparameter evolution with GA	46
4.1	Products in the dataset	49
4.2	Correct vs Incorrect Annotation	53
4.3	Surface textures for synthetic image generation	54
4.4	Synthetic Image	55
4.5	Distribution of products in dataset splits	57
4.6	Distribution of Synthetic and Authentic images in training set	58
4.7	Comparison of Yolov5 Models	60
4.8	Yolov5 Architecture [10]	61
4.9	Comparison of Yolov5 Models	63
4.10	Wireframe of the app	64
4.11	GUI of the Application	65
4.12	CV based billing system	67
4.13	Generated Receipt	68
5.1	Correlogram of the dataset	70
5.2	Intersection over union	71
5.3	Confusion Matrix	72
5.4	Samples from Test Set - 1	75
5.5	Distribution of Test Set -1	76
5.6	Distribution of Test Set -3	77
5.7	Training Statistics for Stacked_Medium model	78
5.8	Confusion Matrices for Front facing products in Test set-1	79
5.9	Confusion Matrices for Back facing products in Test set-1	80
5.10	Comparison of mAP in different lighting conditions	82
5.11	Selected Frames for the Test-2	83
5.12	Analysis for time frame 0:33	83
5.13	Analysis for time frame 0:51	84
5.14	Analysis for time frame 1:15	84
5.15	Analysis for time frame 1:48	85
5.16	Confusion Matrices on Test Set - 3	86
5.17	Comparison of mAP after Data Augmentation	88

5.18	Distribution of Hyperparameter values for 150 generation of evolution	89
5.19	Comparison of mAP after hyperparameter evolution	92
5.20	Frame Analysis for Test Set 2 after Hyperparameter Evolution . .	93

LIST OF TABLES

2.1	Different feature descriptors of different categories used to extract image information	5
3.1	Comparison of speed between R-CNN, Fast R-CNN and Faster R-CNN	25
4.1	Versions of datasets	49
4.2	Pretraining information on the COCO dataset	60
4.3	Hardware used to train the models	62
4.4	Training settings	62
4.5	Shortcuts for the App	66

CHAPTER 1

INTRODUCTION

Supermarkets or departmental stores are an integral part of human life all across the world. A lot of people go there in order to shop for their daily necessities. With a lot of people comes the problem of long queuing times in the checkout corner. Most supermarkets implement traditional infrared barcode-based scanners for billing. One case study of a Brazilian supermarket shows that a customer has to wait for 1.9-1.6 minutes on average in the queue and their checkout process takes around 4.5 minutes of service time from an average cashier.[11] While increasing the number of cashiers can reduce queuing time, the service time is completely dependent on the cashier's speed. This project aims to speed up checkout time by replacing traditional barcode scanners with a computer vision (CV) based checkout solution based on deep learning.

1.1 Motivation

Supermarkets offer a huge array of different variety products. Most stores use traditional billing methods based on infrared barcode scanners. Usage of bar code scanners is heavily time-consuming due to the workflow of using the scanners. The steps a cashier has to perform in order to scan the products are as follows:

- Pick up a product from the checkout area

- Search around the packaging for the barcode
- Scan it with a barcode scanner, if for some reason the scanner isn't immediately scanning the code, then the cashier has to reorient the scanner or the product to aid the scanning.

Picking up and searching for bar codes on each product individually, take up a lot of time. The process could be a lot faster if a batch-billing method is implemented where a batch of products will be billed at once instead of billing them individually. With the rise of deep learning-based object detection algorithms, image processing-based batch-billing of supermarket goods is becoming more and more viable. Companies like Amazon are already implementing such technologies in shopping carts called Dash cart[12]. But adding a computer-based system to each shopping cart is a very expensive endeavour and is not suitable for implementation in developing countries. So in this project, we aim to make a more viable solution by proposing a computer vision-based billing system at checkout corners.

1.2 Objectives

We aim to build a CV-based billing system for supermarkets. It is worth mentioning that the system we propose is not a self-checkout system. A cashier is still required to bill the products, we only aim to replace traditional barcode scanning with a CV-based solution that enables batch-billing. The objectives we aim to achieve are as follows:

- Implementing a CV-based system that can correctly classify supermarket products in batches with a view to reducing checkout time.
- Ensuring lighting independent product detection
- Ensuring that products can be detected both from their front and back labels making the system orientation independent to a certain extent
- Creating a dataset of locally available supermarket products
- Implementing the system in a GUI-based application for ease of use

CHAPTER 2

RELATED WORKS

2.1 ARC Vision-based auto retail

In this study [13], the authors implemented a motor-driven conveyor setup that carries one product at a time inside a wooden hood, which has a webcam placed over its roof. Inside the hood, a Laser LDR (Light Dependent Resistor) module sends signals to a microcontroller board if a product has arrived inside or not. Depending on these signals, the board controls the switching of the conveyor motor. On detecting the product, the microcontroller communicates with the Python environment on the computer to access the webcam to extract a frame and process the image with OpenCV. In the processed image, they used Keras for object identification. They developed a GUI system using Tkinter that facilitates adding the identified products to the shopping cart with their price. In the CNN architecture, they employed ReLUs(Rectified Linear Units) [14] as the activation function, max-pooling, weight initialization [15], dropout [16] and batch normalization. For their training samples, they used 31,000 images of 100 different local retail products, and it produced a training accuracy of 94.76%, validation accuracy of 95.24%, and overall accuracy of 91.7%.



FIGURE 2.1: Structure of the hardware used in the paper

Limitations:

- Product detection one at a time consumes too much time and hence doesn't meet the goal of replacing the current manual retail check-out system.
- The limited size of the wooden hood restricts the size of the products.
- The system is limited to detecting the product in an ideal alignment condition.

2.2 CV based Retail Store Product Detection

The goal of this study [17] is to build a vision system for supermarkets that would enhance the user experience for the consumer and impact on the commercial benefit. This paper works towards an ideal vision system that can produce the availability in a specific product in the inventory, to create a planogram of the product provided the live feed display of the merchandise and better the overall experience of customers that can add value to the commercial benefit. To progress toward meeting these goals, this article researches various methods to enhance product detection systems.

Firstly, the features to be extracted from the images of the products can be achieved using various algorithms. Table 2.1 shows the different categories and descriptors used for the feature extraction from the images of the products:

TABLE 2.1: Different feature descriptors of different categories used to extract image information

Catagories	Feature Description
Key point based Features	SIFT, SURF, BRIGHT
Gradient based Features	Morphological Gradient, HOG, Sobel operator, Canny Edge Detector
Pattern based Features	Haar-like Features, Recurring Patterns
Color based Features	Color Histogram, Saliency, Color Constancy Model
Deep Learning based Features	CaffeNet, AlexNet, Inception-V3, VGG-f, CNN

This paper thoroughly researches different methods for detecting products. Even though detection of products on a daily basis is a simple task for the human eye, the CV faces many challenges due to the recurring shape of identical but different products of a supermarket. Hence the detection method for the products of a supermarket must be tested and studied sincerely. The methods discussed and experimented:

2.2.1 Block based Methods

Block-based method can be divided into two categories, the sliding window method and the grid base method. The sliding window method uses publicly available databases for unsupervised learning. In a recent study [18] Ray et al. approached this dataset with a two-layer method, where the first layer uses SURF [19] to extract the features from the image, the second layer works with a graphical approach to find the best possible match. The grid-based method used both supervised and

unsupervised methods for detection. SIFT [20][21], and BRIGHT [22] are the extraction methods used for unsupervised learning method. For supervised learning method, SIFT, LLC [23], OCR [24] and SVM (used in active learning [25] and used for different attempts.

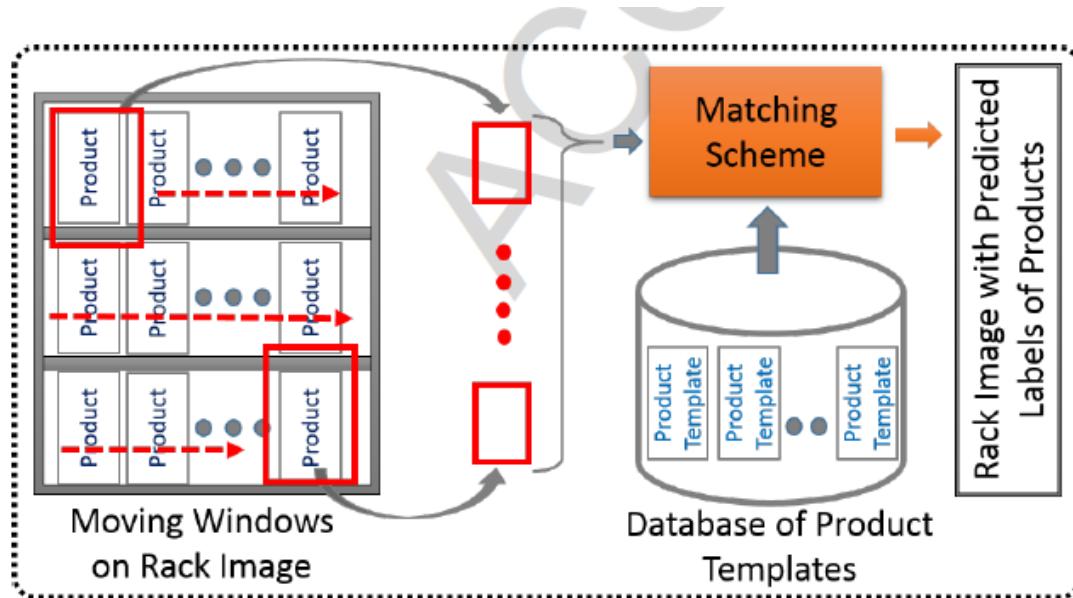


FIGURE 2.2: Block representation of sliding window

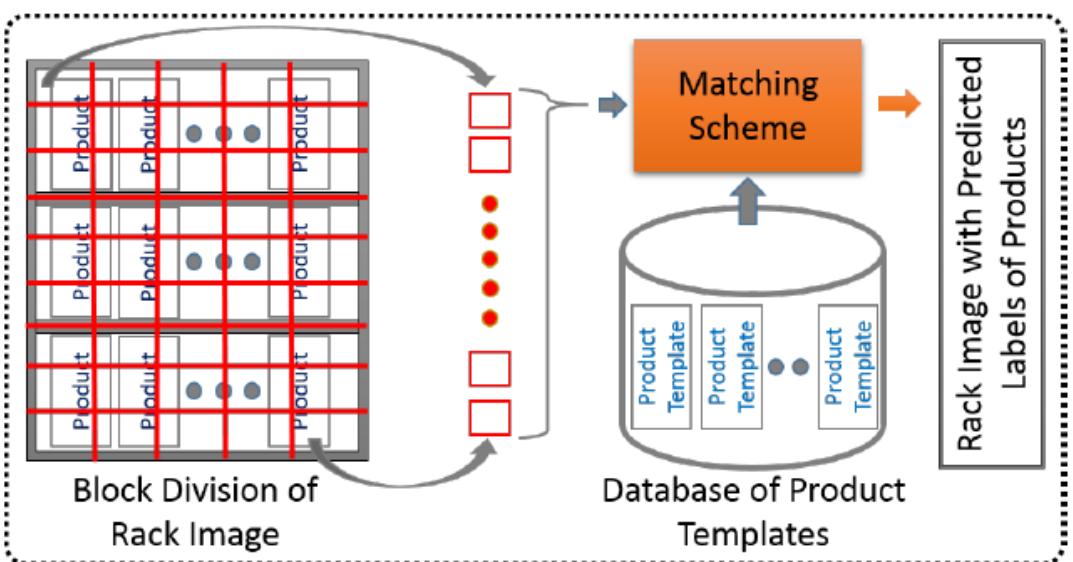


FIGURE 2.3: Block representation of grid based method

2.2.2 Geometric Transformation based Methods

Geometric transformation based methods work toward detecting an image regardless of its angle from the camera device. This method also uses both supervised and unsupervised learning, while using most approaches use the key point-based local features on the databases.

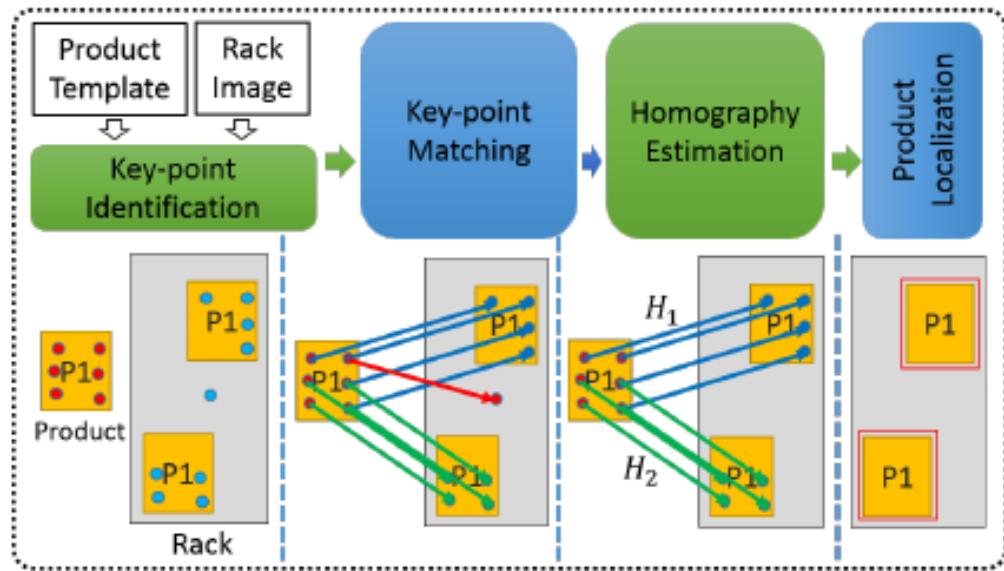


FIGURE 2.4: Block representation of geometric transformation based method

Where P1 is the product image H1, H2 are the homographies and The points represent key points of the image

2.2.3 Saliency based Methods

This method used a number of implementations to detect the product in a shelf of the market, namely saliency maps [26], gradient image [27], potential regions [28] and others. This utilization can also be based on both supervised and unsupervised learning.

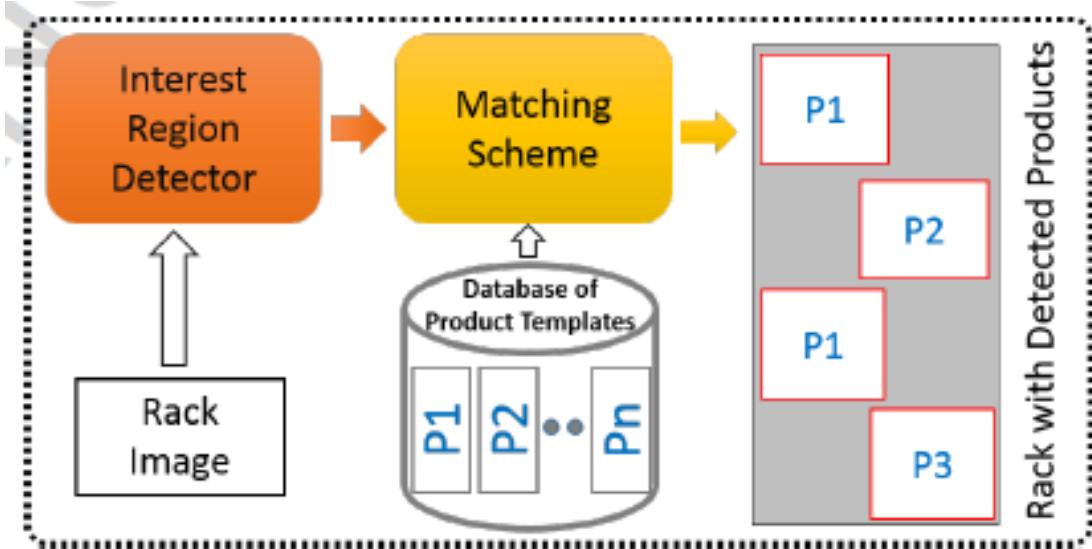


FIGURE 2.5: Block representation of saliency-based method

2.2.4 Detector based Methods

Detector based methods firstly detect the object in a specific image, creating a bounding box around it. A layered extraction is then implemented on the image to extract the features of the image, which then is compared to the product image to classify. A recent work of Karlinsky et al.[29] uses SIFT to extract data from training images, which produces a primary detection. The primary detection is then run through a CNN model (VGG-f network [30] to classify the product. A KLT [31] tracker keeps track of the product detection in a video feed.

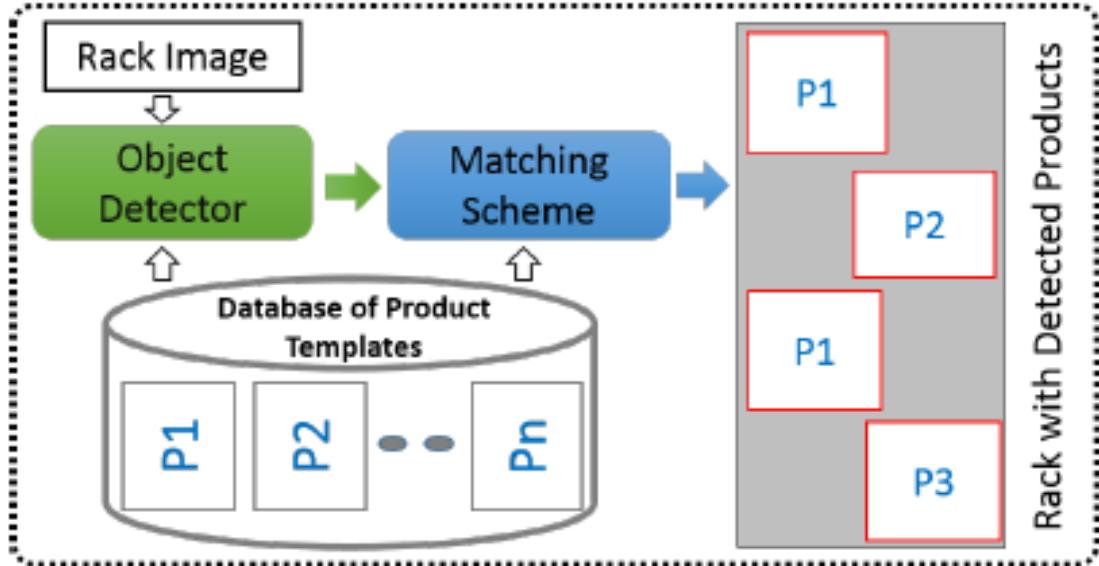


FIGURE 2.6: Block representation of detector-based method

2.2.5 User-in-the-loop Methods

With this method, the bounding boxes of the images are manually added by the user. The image is then cropped either manually or by using a planogram, which then implements local feature extraction on the images and matches it with the product data. This method also works with both unsupervised and supervised data. [17]

2.3 A Computer Vision System Supporting Blind People

This study [32] proposes e-vision utilizing a computer vision application for going to the supermarket to buy food considering handicaps faced by blind people. Despite significant advances, assistive devices based on CV technology for vision-impaired persons are still limited to recognizing barriers and general objects without taking into account the context of the individual's actions. The functional requirements of an assistive device are considerably different in this scenario.

In order to provide a user-friendly experience, they examine the abstraction levels of information that must be communicated in the exhibited system while visiting a supermarket. For example, If the user is looking at a trail, shelf, or product, or if they are at the supermarket's entrance/exit. The user would require various amounts of information at each abstraction level, on the trail information, the system should be able to indicate the specific trail you are on, for example, the drinks trail, added with the information of the level you are looking at i.e. the beer level and should be able to specify you about the brand of the beer you're holding- on the product level.

For the hardware part, this study uses three sensory devices:

1. On the person's head or glasses, or on any body part on the same level of height, a camera is set.
2. To get images a mobile device application is used that the user carries, which gets the information through the designed methodology.
3. To communicate auditory information, a set of earphones is used, with pairs with the mobile application, giving feedback from the camera feed, which is processed through the mobile application, using text-to-speech technology.

The information extracted from the image feed is distinguished in 4 abstract levels, namely Product, Shelf, Trails, Other. The trails give the user the information of what section of products they are in front of. The shelf is information of the user being in front of a specific category of the above-mentioned products. And at the base level, the abstract Product communicates the information on the specific product (i.e. its specific name with its brand).



FIGURE 2.7: Sample images for different levels of abstract

For achieving this critical goal, this article makes use of the improvised Deep Neural Network(DNN) paired with Support Vector Machines (SVMs). Comparing three DNN architectures, VGG16 [33], ResNet [34] and Inception [35], pre-trained on ImageNet Datasets, the features from images were extracted. After categorizing the image on an abstract level, a mechanism is used to extract the accessible texts from the product packages, namely Optical CharacterRecognition (OCR) mechanism which can identify characters both in English and Greek.

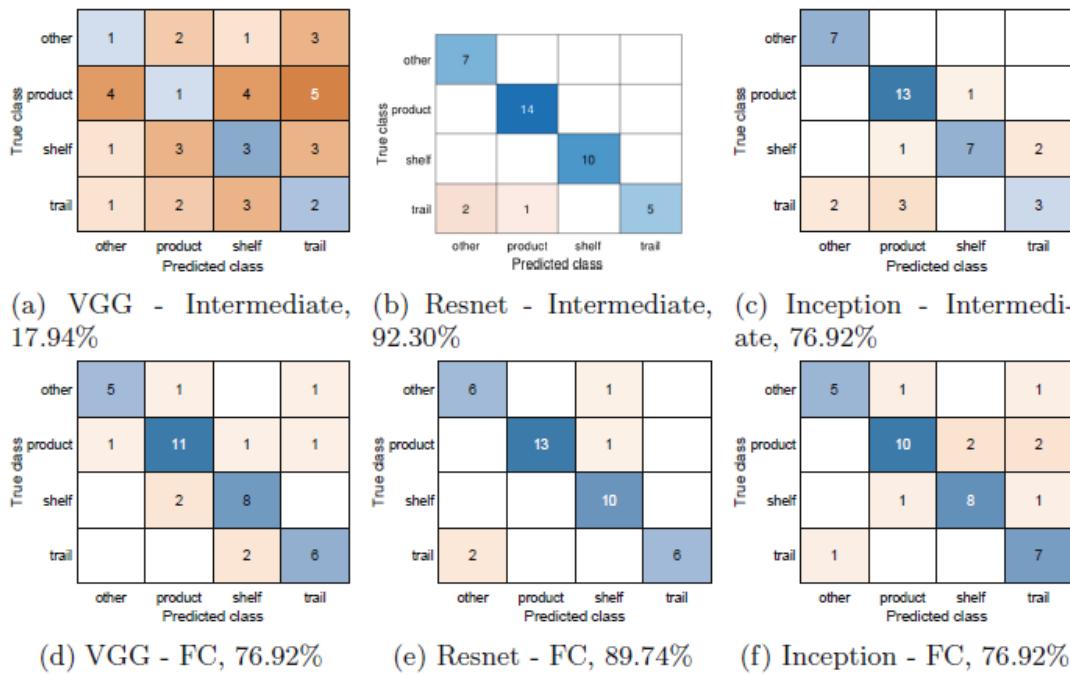


FIGURE 2.8: Representation of Confusion matrices for the quaternary classification task using VGG, Resnet, Inception

This article achieved an accuracy of 77.15% in communicating the information to the user.

Limitations:

- Inadequate dataset for on-site visits that limits the accuracy improvement
- “Other” cases are not defined properly (i.e. for shopping carts or employees)
- A visually assisting system build for ease of payment is not included

- Needs a system for differentiating packaged and non-packaged products

2.4 Toward Real-Time Grocery Detection for the Visually Impaired

The study [36] proposes ShelfScanner, an object detection system for visually impaired customers facilitating them to shop for groceries without asking for another human guidance. This system can detect products on the shelves in real-time from the shopping list of blind people by scanning a store area using only a mobile phone. It implements a translational motion model considering the approximate planar nature of the grocery shelves. Their training data consists of in-vitro images collected from the GroZi-120 database that provides 5.6 images on average per product. Their testing data consists of in-situ images captured from live video of the mobile camera. The system works by detecting a probable set of points from the camera's FOV(Field of vision) when these match any product of the shopping list and notifies the customer. ShelfScanner can also guide the customer for the whereabouts of any previously detected product by implementing a mosaic [37], which is continually being updated by adding new frames data with the running FOV. They did a cycle of the online evaluation for their system, which employs the Lucas-Kanade optical flow method(mosaicing) by OpenCV, SURF for interest point detection [19], NIMBLE [38] for estimation of the probability distribution over the shopping list classes, keypoint selection by comparing it with a threshold. Comparing with a strict threshold, out of 52 item classes, their system identified 17 easy items with zero false positives, eight moderate items with 1-100 false positives and hard items having more than 100 false positives.

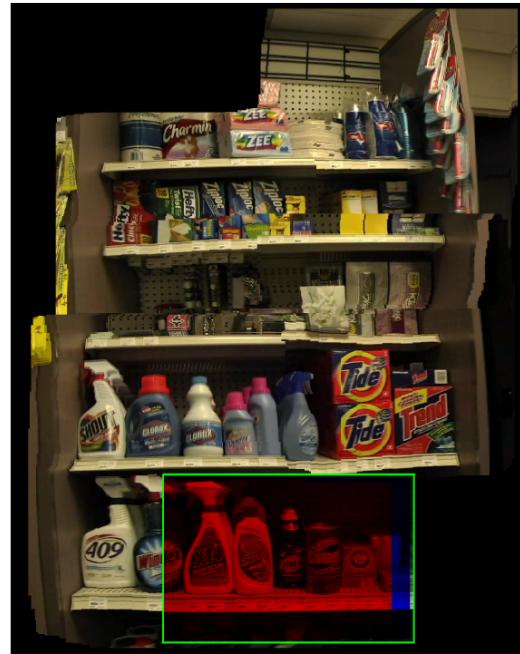


FIGURE 2.9: Under-construction state of the mosaic represented by the green box in the FOV(Field of Vision) of the camera



FIGURE 2.10: Training (top) and testing (bottom) examples of (a) easy retail items, which gives zero false positives for a strict threshold; (b) hard retail items, which fails to detect with reasonable false positives

Limitations:

- The training data are dependent on GroZi-120, which isn't always up to date with the current product appearances.

- The system operates at 2 FPS which is too slow for real-time application.
- As the system maintains a large mosaic for the slow FPS, the memory cost is high.
- Implementation of translational motion model affects object detection.
- The training images are obtained in ideal conditions but not in the real environment, for which the system couldn't detect hard retail items with reasonable false positives.

2.5 Fine-Grained Grocery Product Recognition by One-Shot Learning

This study [39] introduces retail product detection by combining feature matching with single-shot deep learning algorithm using a camera in the grocery stores. They implemented an unsupervised learning algorithm [40] for detecting the recurring features between the captured image of the product on the shelf and the sample image of the single-shot training algorithm. They produced a candidate ROI (Region of Interest) of the product on the shelf that contains only the distinguished features like the logo area and disregard the non-salient (redundant) attributes. They used the RANSAC algorithm [41] to align the feature map of the candidate ROI with the training sample to generate an attention map. The attention map magnifies the fine attributes of the product and feeds the data to VGG-16 (CNN classifier). They collected their dataset from GroZi-3.2 [28], GroZi-120 [42], GP-20, GP-180 [43] and CAPG-GP [39] servers. They tested using SIFT, C-SIFT [44], OpponentSIFT [45], RGB-SIFT, BRISK [46], SURF and found that the model employing SIFT along with VGG-16 showed better balance between precision and recall in detecting the product.

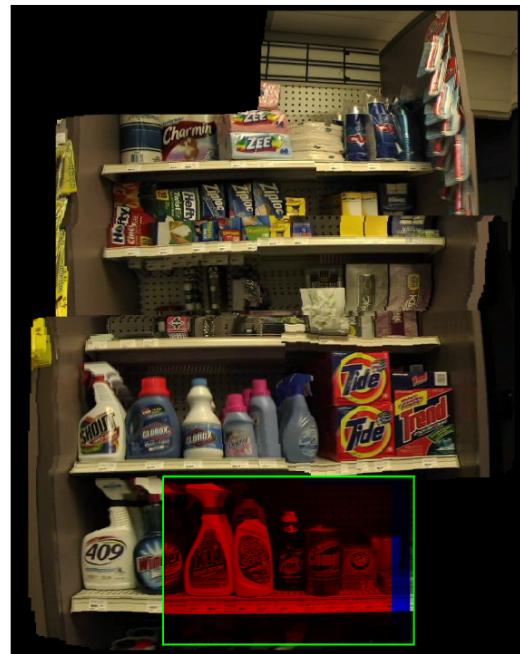


FIGURE 2.11: Proposed framework of the study



FIGURE 2.12: Production of candidate region by extracting feature map of the retail items on the shelf and comparing it with the feature points of the training sample

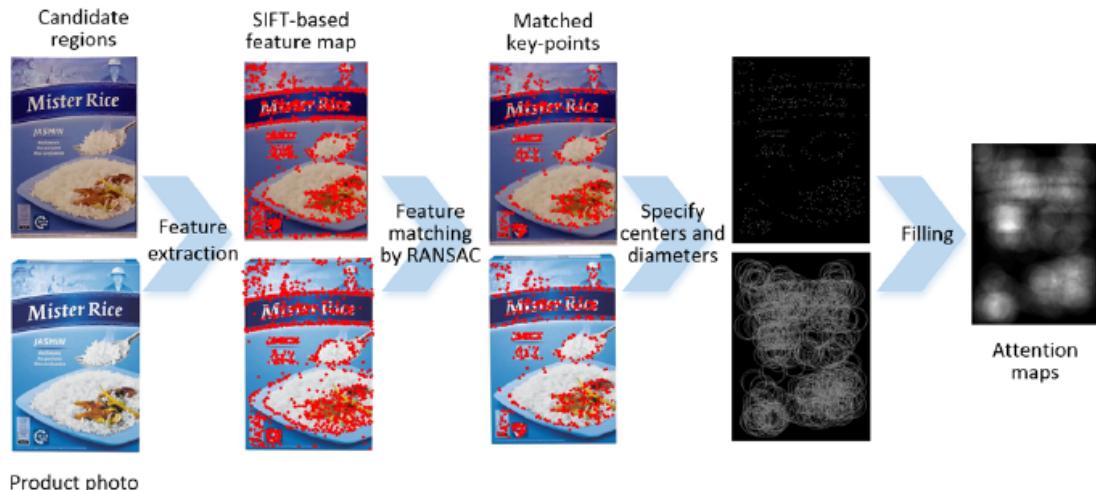


FIGURE 2.13: Generation of Attention Map Using SIFT

Limitations:

- The system fails to detect retail items when not aligned properly in a front-facing way.
- The system does not perform well on detecting the products, which reflect highlights due to the glossy nature of the packaging.

CHAPTER 3

THEORETICAL OVERVIEW

In this chapter, we discuss the necessary theoretical overview of the project along with specifications for any hardware tools used.

3.1 Hardware Specification

We use a Logitech 930e webcam to take the live feed of products on the checkout table. Any other standard USB webcam can be used in its place. The specifications of the webcam are listed below [47]:

- Up to 30fps HD 1080p video quality
- 90 Degree Field of view
- ZEISS glass lens
- Autofocus
- 4x Digital Zoom
- Noise cancelling dual mics
- Dimensions Width 3.7 inches (94mm) Height 1.7 inches(43mm) Depth 2.8 inches (71mm)

- Weight 5.7 ounces (162g)
- USB A Cable length - 5ft (1.5m)

3.2 Deep learning and Neural Networks

We consider the biological brain as the most well structured system for processing different information through our sensory organs like our eyes, ears, skin, nose, tongue etc. The brain is made of neurons and the connections between these neurons help process complicated high-level information in the brain. A neuron by itself may be very simple but a complex structure developed from these simple units are responsible for all of the information processing that happens inside a human or any animal.

In the realm of artificial intelligence, the term "artificial intelligence" refers to a computer Neural networks are used to simulate the incredible structure of a neural system that can learn patterns from prior data and gradually improve at the task it is given [48].

Neural networks can vary in complexity from a simple neuron to complex structure of multiple layers of neurons. A neuron is characterized by inputs outputs and weights and activation functions. The inputs are multiplied by the weights, and the sum is sent through a transfer function, resulting in the neuron's output [49]. This simplest form of ANN(Artificial Neural Network) model is called a perceptron [50].

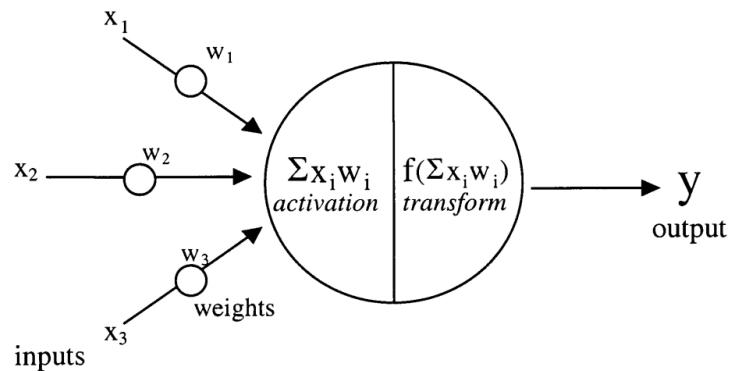


FIGURE 3.1: Model of Perceptron

Complex structures can be created by connecting neurons in multiple layers. So the input and output layers are separated by multiple hidden layers that help learn complex patterns in the data. The gradient of a feed forward NN is evaluated by backpropagation error and is minimized over time as the network iterates.

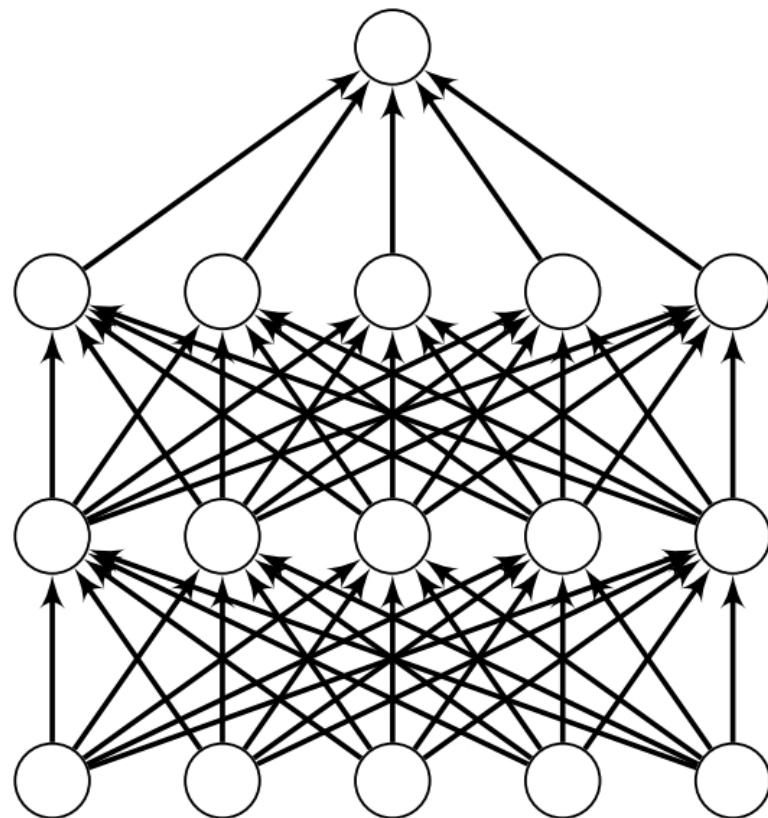


FIGURE 3.2: Standard Deep Neural Network

The error is slowly minimized through use of optimization algorithms like stochastic gradient descent. This leads the model to become more accurate in predicting or classifying the problem set it is designed for [51].

3.3 Transfer Learning

With an abundance of machine learning algorithms focused on training a specific work, it is more efficient to take the knowledge of processes for completing a single task, and use it to learn similar work. This process of using the knowledge of completing a task to help learn to complete a different task is Transfer Learning. Most of the deep learning today is used for completing an isolated work. However, building a new algorithm for a new task is not only time consuming but also inefficient since this process needs to run through a lot of trial and error. By using transfer learning we can accelerate the speed of learning at an efficient rate.

Some of the works for transfer learning are included in neural networks, Bayesian networks, Markov Logic Networks, Q-learning and policy search [52]. While it seems like using training data and future data in the same field of work might be more relative, using transfer learning to train a dataset on a completely different classification task's knowledge can be impressively improving in the context of performances. In real-world problems, algorithms for each different task might not work as well as it would in the case of implementing a transfer learning on that task[1].

Transfer learning can improve the learning process significantly from three contexts:

- The initial learning rate can greatly improve by using transfer learning. While a blind algorithm with no knowledge of completing a task starts off by randomly mathematical approaches, a transfer learning method can help accelerate the slope of initial learning that has already worked on a different task before.

- The saturation point of a learning algorithm reaches much earlier for transfer learning, with a higher performance too. While the amount of time to learn from scratch is significantly longer for learning it from scratch.
- Multi-Task learning is a process where different tasks share knowledge with each other. In this category, information can flow freely between tasks making it easier and faster for them to learn.

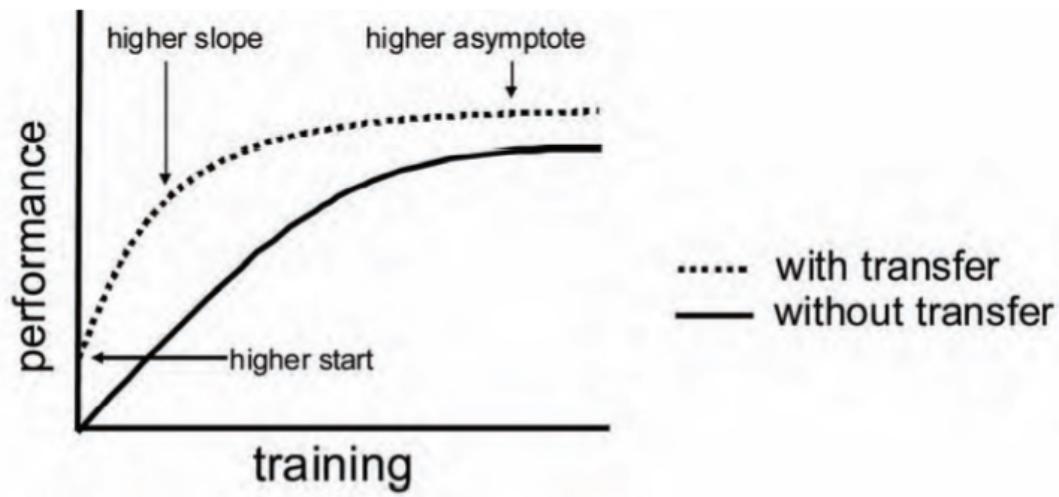


FIGURE 3.3: Improvements due to transfer learning [1]

3.4 Defining Object Detection

In the field of computer vision (CV) the term "Object Recognition" is used to generalize an array of tasks that are related to identifying objects in an image. There are tasks such as Image classification, Object localization and segmentation that fall in this category of tasks.

Object detection has been in the research field for a very long time. The first era of object detection leaned on cold weapon detection, where most of the features were handcrafted. Around 2012, first CNN based object detection was used [53][54] and since then have constantly evolved. Today we use advanced DNN for object detection.

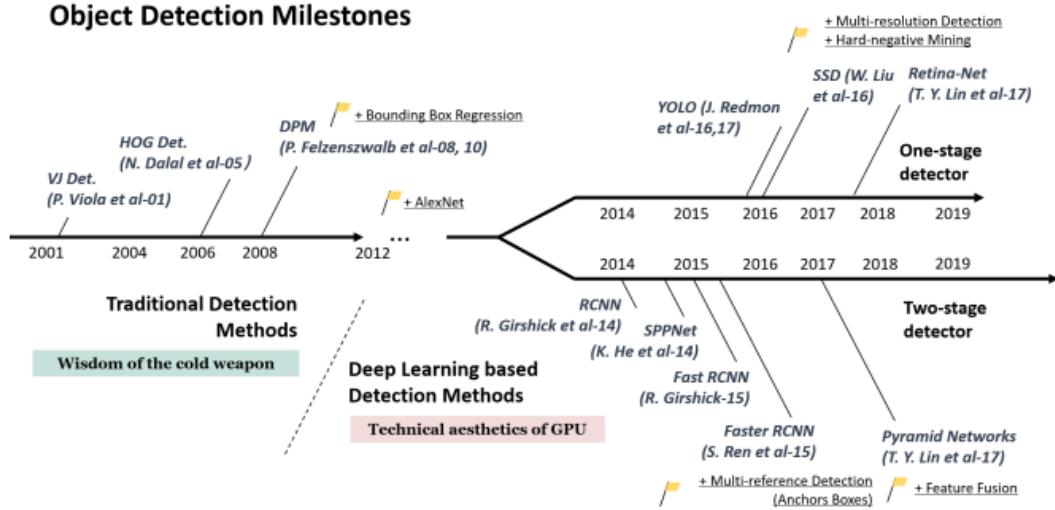


FIGURE 3.4: Roadmap of Object Detection [2]

Using DNN in object detection has shown incredible performance in object detection. For detecting an object we take an image, run it through specified layers of DNN, which then produces masks for the object along with the portions of the object. This process localizes the objects location precisely [55].

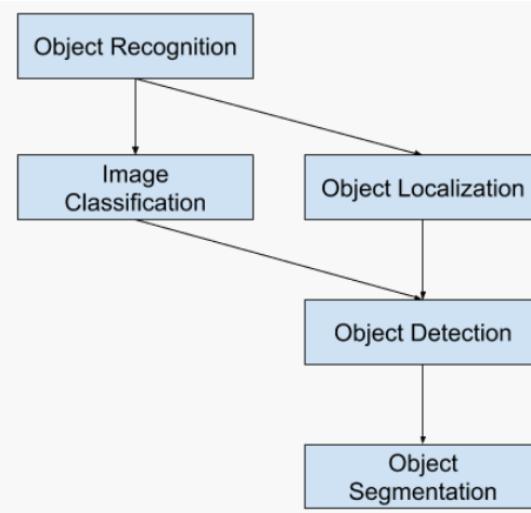


FIGURE 3.5: Core Tasks of object detection

The core tasks of object detection can be portrayed below tasks:

- Image classification refers to the task of identifying what the object in an image is or in other words which predefined class does this object belong to.
- Object Localization is the task of locating where the objects in an image are. This is most commonly done by drawing bounding boxes around them.
- Object Detection is a combination of both these tasks, it first localizes the objects in an image and then identifies which class each object belongs to.
- Object Detection is a combination of both these tasks, it first localizes the objects in an image and then identifies which class each object belongs to.
- Object Segmentation / Semantic Segmentation is an extension of object detection. Instead of drawing bounding boxes around an object, it colours the pixels encompassing each object in a unique colour, visually separating them from each other [56].

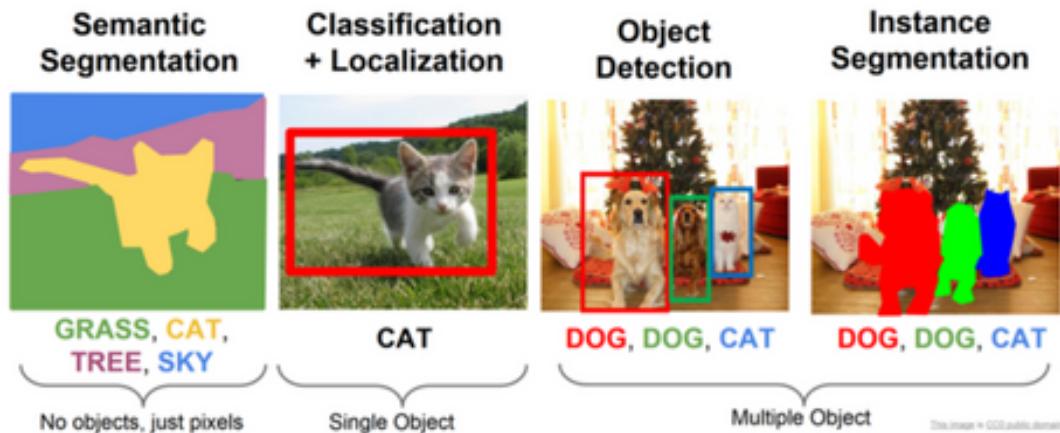


FIGURE 3.6: Comparison of semantic segmentation, classification and localization, object detection and instance segmentation. [3]

3.5 Object Detection Algorithms

3.5.1 R-CNN

R-CNN stands for region-based convolutional neural network which is a neural network specifically used in machine learning for the purpose of object tracking. Computer vision has seen rapid growth in technological aspects largely due to baseline systems such as fast/faster R-CNN [57]. It is known to provide state-of-the-art results in the field of object tracking. Mask R-CNN architecture also exists as an intuitive extension of fast R-CNN which adds a branch for each Region of Interest, forecasting segmentation masks [58].

The primary purpose of R-CNN is to identify the primary objects of any given picture. The detection process is made up of three different modules. The first creates proposals based on independent regions. The second module is a feature extractor implemented through deep CNN. The third module is for classification, and it may follow any classification algorithm such as kNN, SVM etc. [59]. The R-CNN pipeline proposes a set number of boxes per picture that are most likely to contain the target items. All proposal boxes (even small ones) are enlarged to a canonical size at the last pooling layer, which means that a complete feature map is created for each proposal box [60].

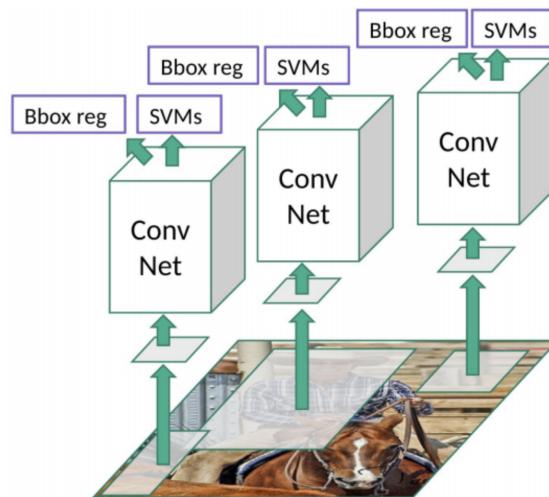


FIGURE 3.7: R-CNN Structure RCNN

The advantage of R-CNN is that it is highly accurate, but it is also prolonged and unsuitable for real-time applications. The difference in test times is shown in the following table. Fast and Faster R-CNN works much better for applications where real-time input is involved.

TABLE 3.1: Comparison of speed between R-CNN, Fast R-CNN and Faster R-CNN

	R-CNN	FastR-CNN	FasterR-CNN
Test time perimage(with Proposal)	50 sec	2 sec	0.2 sec
(Speedup)	1x	25x	250x
mAP(VOC 2007)	66.0	66.9	66.9

3.5.2 SPP Net

The SPPnet technique creates a convolutional feature map for the whole input picture and then uses a feature vector taken from the shared feature map to classify each object proposition [61]. At test time, SPPnet speeds R-CNN by 10 to 100 times. Due to quicker proposal feature extraction, training time is also decreased by three times. However, SPPNET exhibits similar problems as R-CNN. Extracting features, fine tuning a network using log loss, training SVMs, and ultimately fitting bounding-box regressors are all part of the training workflow. However, it is not similar in that the convolutional layers that precede the spatial pyramid pooling cannot be updated using the finetuning technique. This limits the precision of very deep neural networks [62].

3.5.3 Fast RCNN

Although RCNN is highly popular for object tracking, it is still far too slow for real-time applications. Fast RCNN is known to have much better performance with minimal delay. The primary issues of using RCNN are:

1. R-CNN uses log loss to finetune a ConvNet on object suggestions. The SVMs are then fitted to ConvNet features. These act as object detectors,

and they replace the classifier learnt by fine tuning. Finally, in the last stage, regression of the bounding boxes is discovered. This means there are multiple stage pipelines to go through before the data can be processed.

2. There are multiple object proposals based on a region, and each has to be extracted in order for R-CNN to work. Considering the high dimension of modern data it would take a very long time and a lot of GPU memory to allow the CNN to work properly.
3. Detection of the average of a VGG16 image takes 47 seconds. This means the detection process is very slow.

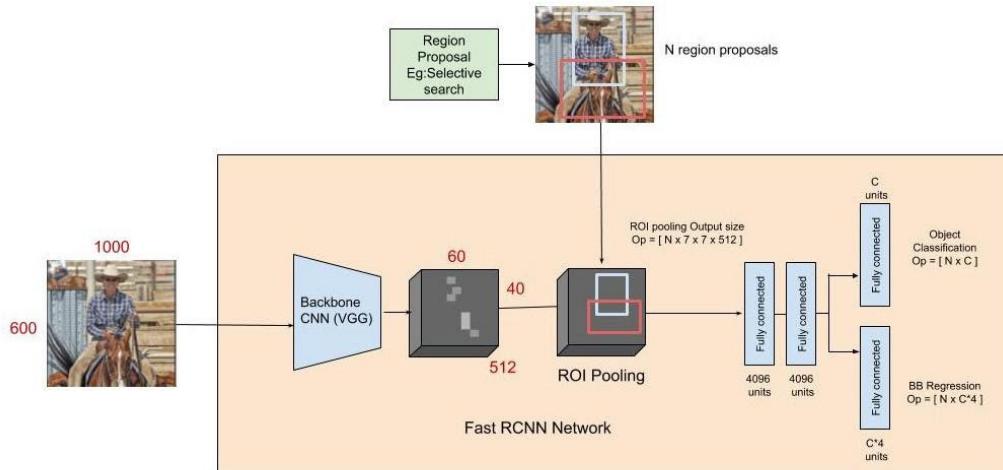


FIGURE 3.8: Faster RCNN architecture

A full picture and a collection of item suggestions are fed into a Fast R-CNN network. To create a conv feature map, the network first analyzes the entire image with multiple convolutional (conv) and max-pooling layers. After that, a region of interest (RoI) pooling layer generates a fixed-length feature vector from the feature map for each object proposition.

3.6 YOLO Overview

3.6.1 History

The YOLO (You Only Look Once) algorithm was introduced in 2015 by researcher Joseph Redmon and others using a novel approach, treating object detection as a regression problem using only a single neural network. It brought remarkable achievements in the field of object detection than other models could. Five different versions of YOLO have been developed hitherto, of which the first three were developed by the original authors, Joseph Redmond and others. However, he discontinued his research in the computer vision field after the release of version three for ethical reasons. In 2020, Alexey Bochkovskiy, who developed the previous three versions of YOLO with Redmond, released version 4, YOLOv4 [63], on the official YOLO Github account. Glenn Jocher and his Ultralytics LLC research team developed the YOLO algorithm using the Pytorch framework and released the YOLOv5 only after one month of the release of YOLOv4 with a few changes and improvements. Although none of the original authors was involved in developing this version, the YOLOv5 outperforms all the previous versions in terms of both accuracy and speed. Since no paper exists yet of YOLOv5, there is controversy in the research community if this model justifies the original branding of YOLO [64].

Before YOLO, the other CNN models like R-CNN employed RPPNS(Regional Proposals Networks), which was an iterative method for classifying different regions of the input image. For each region, the model produced a bounding box, ran a classifier on the box, applied post-processing and refined the bounding box. Optimizing detection performance on individual stages of the networks separately made it very difficult [65].

YOLO Redmon_2016_CVPR [4] came up with an approach to unify all stages on a single neural network. The entire detection pipeline of YOLO consists of a single network facilitating the network to be optimized for detection performance easily from end to end. By anticipating the bounding boxes and class probabilities, the network directly analyzes the image and recognizes not only what items are present but also where they are located inside the image. Rather than using an iterative

method, the system can make predictions for all of the objects in the image at once, hence the name YOLO (You Only Look Once). Even though it sometimes struggles with smaller objects, YOLO being orders of magnitude faster than other object detection algorithms mentioned above is our reason for choosing it for the project. [66]

3.6.2 How it works

YOLOv1 applies grid cells of size $S \times S$ (7×7 by default) on an image. The grid cell in which the object's center is located is responsible for detecting that object. Despite the fact that the object appears on other grid cells, those cells are discarded as the center isn't present inside them.

Each grid cell predicts B bounding boxes, each of which consists of some parameters (center, width, height) and confidence score. The confidence score is used to determine if an object is present or absent in that specific bounding box. The confidence score can be calculated as:

$$\text{confidence score} = p(\text{Object}) * \text{IOU}_{\text{truth}}^{\text{pred}} \quad (3.1)$$

Where $p(\text{Object})$ is the probability of the presence of an object inside the cell and $\text{IOU}_{\text{pred}}^{\text{truth}}$ is the intersection over union of the predicted bounding box and the ground truth bounding box. If there is no object in a cell, then $p(\text{Object}) \rightarrow 0$ hence the confidence score is also close to 0. Otherwise, the score is close to $\text{IOU}_{\text{pred}}^{\text{truth}}$.

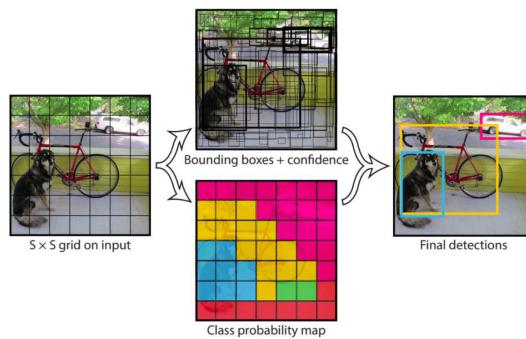


FIGURE 3.9: YOLO model with 7×7 grid was applied to input Image(Redmon, et al. 2016) [4]

Along with the confidence score, each bounding box employs 4 other parameters, which are x, y, w and h. The center coordinate is represented by the x,y parameters, width by the w parameter and height by the h parameter. Therefore, each bounding box uses a total of 5 parameters: x, y, w, h and the confidence score.

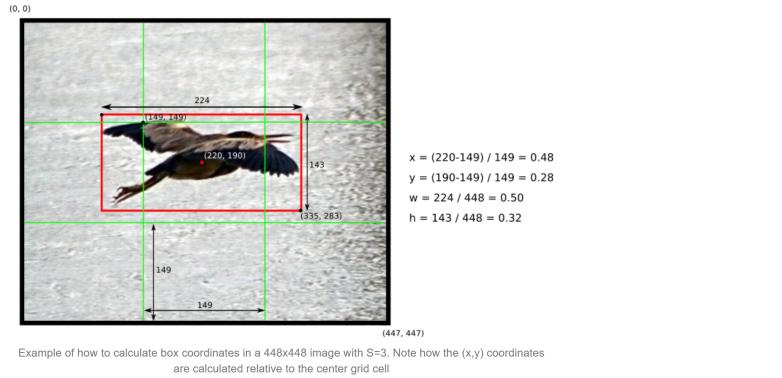


FIGURE 3.10: Example of how to calculate box coordinates in a 448*448 image with S=3. [5]

To predict in which class the object belongs to, conditional probability is applied on the grid cell given that the cell contains one object, $\Pr(\text{Class}(i) \mid \text{Object})$. Therefore, if the grid cell contains no object, the loss function will not make an incorrect class prediction. Regardless of the number of bounding boxes B, the network only predicts one set of class probabilities per cell. In total, there are $S \times S \times C$ class probabilities and an $S \times S \times (B * 5 + C)$ tensor as output which is obtained by adding the class predictions to the output vector.

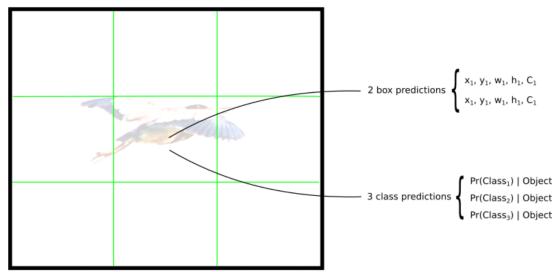


FIGURE 3.11: Each grid cell makes B bounding box predictions and C class predictions. [5]

The predicted bounding box vectors are denoted by output vector \hat{y} and ground truth bounding box vectors are denoted by vector label y . Vector label y and

predicted vector \hat{y} could be indicated in the figure as shown below. There is no object present in the purple cell, hence the confidence score of bounding boxes in this cell is equal to 0, which leads to discarding all the remaining parameters.

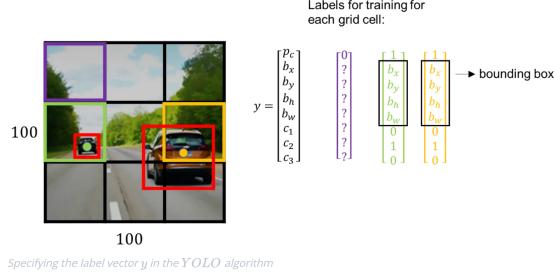


FIGURE 3.12: Matrix calculation for YOLO [6]

Lastly, YOLO removes all the bounding boxes which do not contain any object like the purple cell in the above figure, or contain the same object as other bounding boxes by using Non-Maximum Suppression (NMS). By setting a threshold value, NMS removes all the overlapping bounding boxes which have intersection over union (IOU) value higher than the threshold value [67].

3.6.3 YOLO Overview

YOLOv1 Architecture

The YOLO model was created by the original authors to include Darknet architecture, which analyses all image features, followed by two fully connected layers that perform bounding box prediction for objects. The authors used $S=7$, $B=2$, and $C=20$ in the Pascal VOC dataset to test this model for which the final feature maps are 7×7 , and the output size was $(7 \times 7 \times (2 \cdot 5 + 20))$ [68].

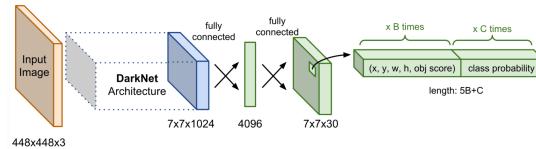


FIGURE 3.13: Network architecture of YOLO [7]

For uncomplicated datasets, the authors introduced the Fast-YOLO model with 9 CNN layers in Darknet architecture, while the normal-YOLO model with 24 CNN layers in Darknet design can handle more complex datasets and produce higher accuracy. Instead of Leaky Rectified Linear Unit (leaky ReLU) activation, the final layer utilizes a Linear activation function [69].

$$\phi(x) = \begin{cases} x, & \text{if } x > 0 \\ 0.1x, & \text{otherwise} \end{cases} \quad (3.2)$$

Name	Filters	Output Dimension
Conv 1	7 x 7 x 64, stride=2	224 x 224 x 64
Max Pool 1	2 x 2, stride=2	112 x 112 x 64
Conv 2	3 x 3 x 192	112 x 112 x 192
Max Pool 2	2 x 2, stride=2	56 x 56 x 192
Conv 3	1 x 1 x 128	56 x 56 x 128
Conv 4	3 x 3 x 256	56 x 56 x 256
Conv 5	1 x 1 x 256	56 x 56 x 256
Conv 6	1 x 1 x 512	56 x 56 x 512
Max Pool 3	2 x 2, stride=2	28 x 28 x 512
Conv 7	1 x 1 x 256	28 x 28 x 256
Conv 8	3 x 3 x 512	28 x 28 x 512
Conv 9	1 x 1 x 256	28 x 28 x 256
Conv 10	3 x 3 x 512	28 x 28 x 512
Conv 11	1 x 1 x 256	28 x 28 x 256
Conv 12	3 x 3 x 512	28 x 28 x 512
Conv 13	1 x 1 x 256	28 x 28 x 256
Conv 14	3 x 3 x 512	28 x 28 x 512
Conv 15	1 x 1 x 512	28 x 28 x 512
Conv 16	3 x 3 x 1024	28 x 28 x 1024
Max Pool 4	2 x 2, stride=2	14 x 14 x 1024
Conv 17	1 x 1 x 512	14 x 14 x 512
Conv 18	3 x 3 x 1024	14 x 14 x 1024
Conv 19	1 x 1 x 512	14 x 14 x 512
Conv 20	3 x 3 x 1024	14 x 14 x 1024
Conv 21	3 x 3 x 1024	14 x 14 x 1024
Conv 22	3 x 3 x 1024, stride=2	7 x 7 x 1024
Conv 23	3 x 3 x 1024	7 x 7 x 1024
Conv 24	3 x 3 x 1024	7 x 7 x 1024
FC 1	-	4096
FC 2	-	7 x 7 x 30 (1470)

FIGURE 3.14: YOLOv1 model architecture. [5]

Loss function

The backbone of YOLO's loss function is the sum-squared error. Multiple grid cells containing no objects have a confidence score of zero, for which the gradients of cells that contain the objects gets overwhelmed. To reduce training divergence and model instability, YOLO applies the largest penalty for predictions from bounding boxes containing objects ($\lambda_{coord}=5$) and the lowest penalty for predictions from bounding boxes with no object ($\lambda_{noobj}=0.5$) [70]. The loss function of YOLO is calculated by adding the loss functions of all bounding box parameters:

$$\begin{aligned} \mathcal{L} = & \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{I}_{ij}^{obj} [(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2] \\ & + \lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{I}_{ij}^{obj} \left[(\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right] \\ & + \sum_{i=0}^{S^2} \sum_{j=0}^B \left\| \mathbb{I}_{ij}^{obj} (c_i - \hat{C}_i)^2 + \lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^B \left\| \mathbb{I}_{ij}^{noobj} (C_i - \hat{C}_i)^2 \right\|_i \right\|_j \\ & + \sum_{i=0}^{S^2} \mathbb{I}_i^{obj} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned} \quad (3.3)$$

Batch Normalization

YOLOv2 implemented batch normalization as a normalization method that facilitates faster and more stable deep neural network training by stabilizing the input layer distribution [71]. This method normalizes the output of each layer after activation to a zero-mean state with a standard deviation of 1.

$$\text{Mini - batch mean: } \mu = \frac{1}{m} \sum_{i=1}^m z^{(i)}$$

$$\text{Mini - batch variance: } \sigma^2 = \frac{1}{m} \sum_{i=1}^m (z^{(i)} - \mu)^2$$

$$\text{Normalize: } z_{\text{norm}}^{(i)} = \frac{Z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\text{Scale and shift: } \tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

This technique decreases training time while simultaneously improving to regularize the network. The network also does not need to use dropout anymore to avoid

overfitting. Batch normalization improved mAP (mean average precision) by more than 2% in YOLOv2 [72].

High Resolution Classifier

In YOLOv1, the model trains the feature extractor(classifier network) at 224×224 and enhances the detection resolution to 448×448 causing the network to simultaneously switch between learning object detection and adapt the new input resolution. Whereas in YOLOv2, the feature extractor is first processed on ImageNet [73] at the full 448×448 resolution for 10 epochs, which allows the model extra time to modify its filters to perform better on higher resolution input than YOLOv1. Then the resulting network is trained on detection which increases mAP (Mean Average Precision) by almost 4% .

Bigger network with ResNet

YOLOv2 failed to detect small objects because the input image lost its fine-grained details due to the downsampling done before entering the deeper layers. ResNet (Residual networks) allowed the activations to enter deeper layers without losing details of input by implementing the concept of skip connections [74].

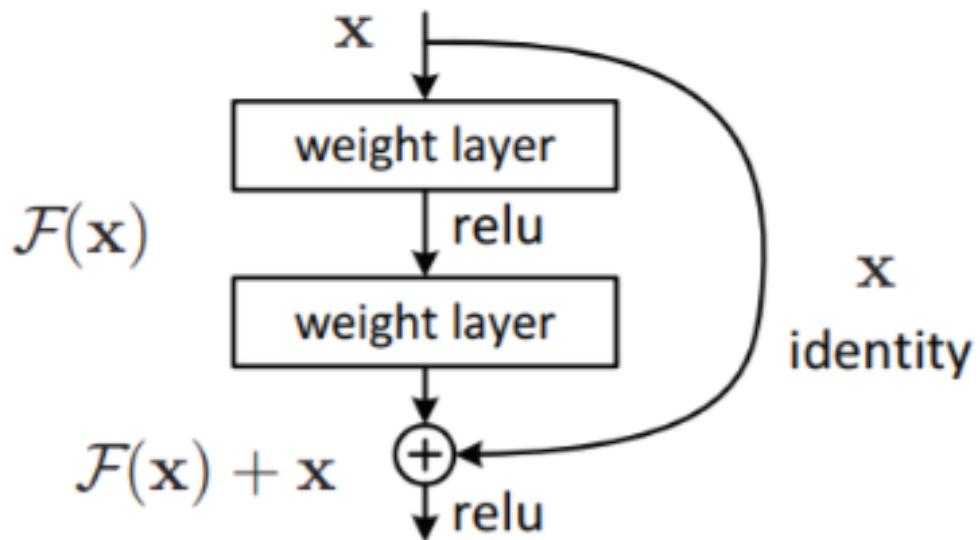


FIGURE 3.15: Building block of residual learning.

Feature Extractor The architecture of YOLOv3 proposes a hybrid feature extractor combining YOLOv2, Darknet-53 and ResNet [75]

Type	Filters	Size	Output
Convolutional	32	3×3	256×256
Convolutional	64	$3 \times 3 / 2$	128×128
1x	Convolutional	32	1×1
	Convolutional	64	3×3
	Residual		128×128
	Convolutional	128	$3 \times 3 / 2$
			64×64
2x	Convolutional	64	1×1
	Convolutional	128	3×3
	Residual		64×64
	Convolutional	256	$3 \times 3 / 2$
8x	Convolutional	128	1×1
	Convolutional	256	3×3
	Residual		32×32
	Convolutional	512	$3 \times 3 / 2$
			16×16
8x	Convolutional	256	1×1
	Convolutional	512	3×3
	Residual		16×16
	Convolutional	1024	$3 \times 3 / 2$
4x	Convolutional	512	1×1
	Convolutional	1024	3×3
	Residual		8×8
	Avgpool		Global
	Connected		1000
	Softmax		

Table 1. Darknet-53.

FIGURE 3.16: Darknet-53.

3.6.4 YOLOv5

The previous versions of YOLO were developed on the Darknet framework, which is written in the C language. However, YOLOv5 is developed in PyTorch, which is written in the Python language. As YOLOv4 and YOLOv5 are developed using

two different languages on two different frameworks, their performance cannot be compared with accuracy. In certain aspects, YOLOv5 has shown better performance than YOLOv4 and has partly gained recognition in the research community besides the other versions of YOLO. Since YOLOv5 is written in the Python language, many applications will be easier to implement.

The YOLOv5 model can be summarized as follows:

Backbone: Focus structure, CSP network

Neck: SPP block, PANet

Head: YOLOv3 head implementing GIoU-loss

CSP Backbone

Both YOLOv4 and YOLOv5 have employed the CSP Bottleneck [76] in their architecture. The CSP addresses the problem of the duplicate gradient of other convolutional networks' backbones, which is crucial for the YOLO architecture as inference speed is extremely important for the YOLO model.

The CSP networks (Cross Stage Partial Network) are based on DenseNet [77]. DenseNet facilitates solving the vanishing gradient problem and reducing the network parameters enabling the YOLOv5 architecture to have a small model size.

Adaptive anchor boxes

The authors of YOLOv2 proposed the concept of anchor box as well as a method for selecting anchor boxes that are similar in size and shape to the ground truth bounding boxes in the training samples. They chose the 5 best-fit anchor boxes for the COCO dataset of 80 classes using the k-mean clustering algorithm and set them as default. However, when a new dataset emerges with a class that is not one of the 80 classes in the COCO dataset, the default anchor boxes are unable to immediately adjust to the ground truth bounding boxes. A giraffe, for instance, requires a thinner and higher anchor box than a square box. Therefore, the k-mean clustering algorithm is used first to get the best-fit anchor box for the giraffe dataset.

In YOLOv5, the anchor box selection process has been integrated. Therefore, this network does not need to consider any of the datasets to be used as input. It automatically learns the best-fit anchor boxes for any type of dataset and uses them during training.[78]

Easy to Use

YOLOv5 is easy to use [79] for computer vision-related applications compared to other object detection architectures since it is written in Python language. Only Pytorch and some lightweight python libraries need to be installed to implement this model. This model can be trained fast, which reduces experimentation costs for the developers.

3.7 Data Augmentation

Deep learning models need a lot of data to learn properly and the acquisition of such large volumes of data are not always possible. With smaller datasets probability of overfitting increases and if we look at the most effective models, we can see that they are clearly driven by the largest datasets. The lack of good quality data can be circumvented by using data augmentation in order to create new data [80]. There are several ways for image based data augmentations,

3.7.1 Basic Augmentation Techniques

Geometric Augmentations:

This includes easy to perform actions like rotation, flipping, translation, random cropping. These methods might be easy to perform, but they sometimes are not label-preserving in nature. So the images have to be relabeled post augmentation.

Colour based Transformation:

Changing the colour spaces for images by manipulating its colour channels in an augmentation technique employed to tackle the lighting based biases in most image recognition problems [81].

3.7.2 Synthetic Data

One of the most important issues in computer vision problems is generating a labelled dataset. Most of the time esp. in case of segmentation tasks the data has to be manually labeled in order to produce good quality results. It is a very tedious process to acquire and label a huge amount of image data manually, so an alternative solution is generating synthetic data [9]. There are many methods for synthetic generation:

Mixing Images:

Randomly mixing different parts of images or changing background based on a set of objects and background images are very simple ways of synthetic data. Images can be randomly cropped and added to other images.

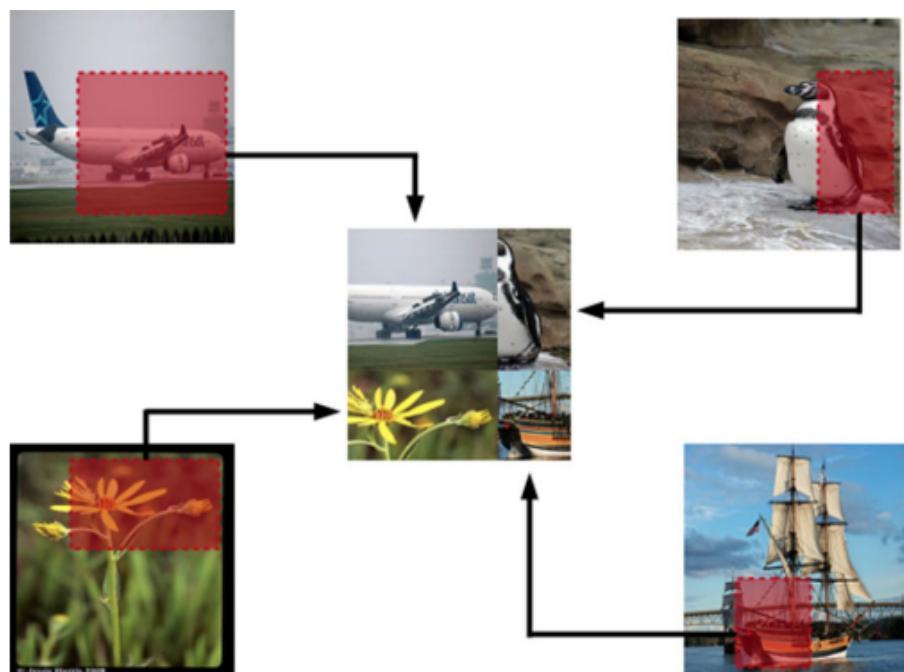


FIGURE 3.17: Image mixing via random cropping.[8]

Another technique for synthetic data generation is composing images by masking objects on different backgrounds. Geometric augmentation techniques can also be applied alongside this method to generate more effective synthetic data [82].



FIGURE 3.18: Image composition via masking

Generative Adversarial Networks(GAN):

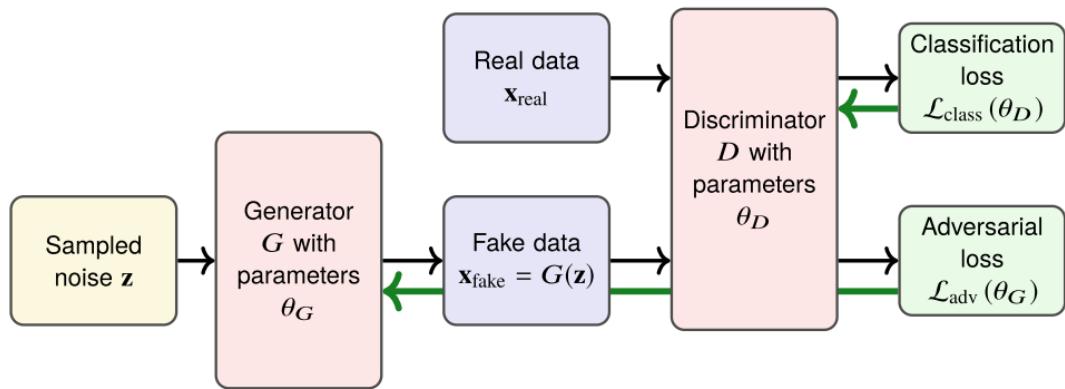


FIGURE 3.19: Basic GAN Architecture [9]

GAN is a framework of two models in an adversarial setting where one model tries to create an image of a sample similar to something from the training set but another discriminative model tries to predict if the generated sample came from the training set or not. Both of the models are trained together with the purpose of making the generative network create entirely new samples of an object from the training data that are significantly different.



FIGURE 3.20: Synthetic image generation with GAN (Yellow marked column is the closest training sample to the neighbouring output.)

3.8 YOLO Annotation

For Object detection algorithms, annotation means identifying the location and class of objects present in an image. This can be done in several different ways, such as bounding boxes, polygonal segmentation, semantic segmentation, 3d cuboids, etc. Bounding boxes are the most common type of annotation used in Computer Vision. An object in an image is identified by the coordinates of a box encompassing its visible area [83].

YOLO has a specific format for annotations. . . YOLO uses the following way to annotate objects in an image. Each object’s bounding box must be represented with its class id, center coordinates, height, and width. The coordinates and dimensions must be normalized by the dimensions of the image. We create a .txt

file with the same name as the image file for each image with the contents like this:

```
cls_id norm_x_center norm_y_center norm_width norm_height
```

$$\text{norm_x_center} = \frac{x - \text{center}}{\text{img_width}}$$

$$\text{norm_y_center} = \frac{y - \text{center}}{\text{img_height}}$$

$$\text{norm_width} = \frac{\text{width}}{\text{img_width}}$$

$$\text{norm_height} = \frac{\text{height}}{\text{img_height}}$$



FIGURE 3.21: Bounding Box

In the Figure we can see the center of coordinates of the bounding box (198.98, 204.01) and the height and width of the bounding box are 120.64px and 333.84px. The image is 416x416. In YOLO's annotation format we need to normalize this values with respect to the dimensions of the image.

$$\text{norm_x_center} = \frac{198.98}{416} = 0.4639$$

$$\text{norm_y_center} = \frac{204.01}{416} = 0.4904$$

$$\text{norm} - \text{width} = \frac{120.64}{416} = 0.29$$

$$\text{norm} - \text{height} = \frac{333.84}{416} = 0.8025$$

Let's assume that the class for this product is 1, So the text file of this image will contain: 1 0.4639 0.4904 0.29 0.8025

3.9 Optimizers

Optimizers are algorithms that work on hyperparameters, the values of which determine the behaviour of the optimizer. In a network, the optimizer determines its learning slope for achieving the intended value of the loss function.

For a loss function produced over a neural network l , its first derivative can be taken as a gradient (theta) can be represented as the hyperparameters. For the above scenario, an optimizer is an algorithm that iterates to reach a point for which l will be minimized. For different optimizer algorithms, the loss function and hyperparameters can differ. [84]

3.9.1 SGD (Stochastic gradient descent)

SGD is a fairly simple but efficient optimizer. The mathematical idea behind this optimizer is to follow the steepest slope to reach the localized minima.

For a neural network that needs to produce a value for the hyperparameter which reduces the loss function to a minimum, SGD approaches the algorithm with below relation:

$$w_k = w_{k-1} - \alpha_{k-1} \hat{\nabla} f(w_{k-1}) \quad (3.4)$$

Where, w_k denotes the k^{th} iterate, α_k is a(tuned) step size sequence, also called as the learning rate.

With f being the loss function. SGDM, a variant of SGD, uses inertia to calculate the learning rate of the optimizer. This can greatly affect the performance of the optimizer. SVGM uses the following iteration to reach a minimum loss point:

$$\begin{aligned} v_k &= \beta v_{k-1} + \hat{\nabla} f(w_{k-1}) \\ w_k &= w_{k-1} - \alpha_{k-1} v_k \end{aligned} \quad (3.5)$$

where

$$\beta \in [0, 1) \quad (3.6)$$

is a momentum parameter and v_0 is initialized to 0. [85]

However SGD may still reach a local minima or a saddle point and its learning rate must be set in between a range to get the minima, or divergence might occur. [86]

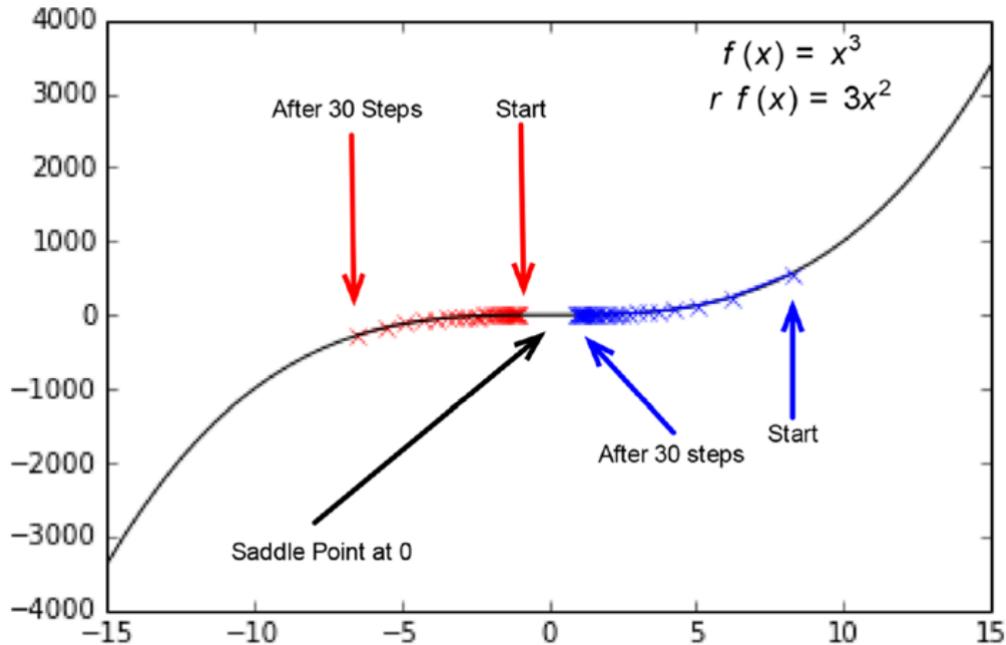


FIGURE 3.22: Local Minima for SVG algorithm

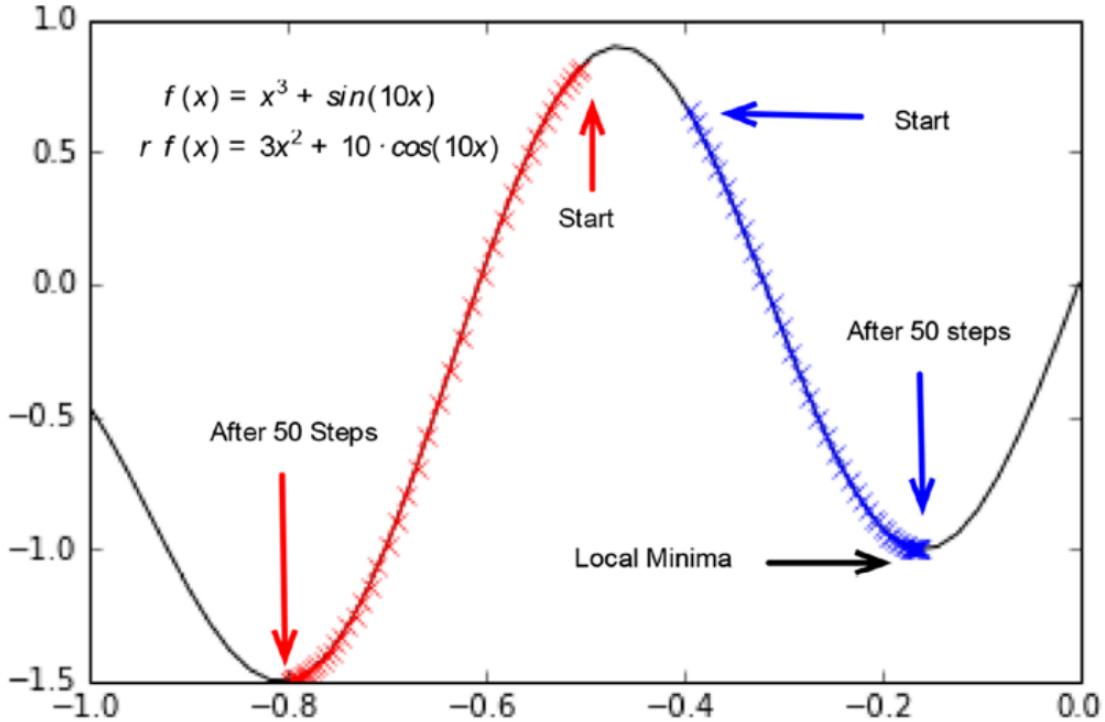


FIGURE 3.23: Saddle Point in minimizing the loss

3.9.2 ADAM

With an attempt to correct these shortcomings, several adaptive optimizers were developed, namely ADAM, AdaGrad, RMSProp. With AdaGrad and RMSProp, the initialization makes the earlier updates messy and can tilt the scaling parameters. ADAM tried to solve this problem using below iteration: $w_k = w_{k-1} - \alpha_{k-1} \cdot \frac{\sqrt{1-\beta_2^k}}{1-\beta_1^k} \cdot \frac{m_{k-1}}{\sqrt{v_{k-1}+\epsilon}}$, where

$$\begin{aligned} m_{k-1} &= \beta_1 m_{k-2} + (1 - \beta_1) \hat{\nabla} f(w_{k-1}) \\ v_{k-1} &= \beta_2 v_{k-2} + (1 - \beta_2) \hat{\nabla} f(w_{k-1})^2 \end{aligned}$$

[85]

In recent studies, ADAM has shown exceptional performance with its adaptive ability. With the generalization and improved training performances, ADAM has solved the problem of inability of SGD to adapt to different scales. [83]

3.10 YOLO Hyperparameters

Hyperparameters are parameters of a machine learning model that are external and cannot be calculated from the data it is operating on. A model tries to minimize some sort of loss function. This process involves some parameters to be kept fixed during the process. Changing these parameters will lead to changes in the behaviour of a model. These parameters, which have to be manually adjusted by humans, are called hyperparameters [87] [88].

YOLOv5 has several hyperparameters; they are as follows [89]:

- Initial learning Rate (lr0)
- Final onecycle LR learning rate (lrf)
- SGD momentum (momentum)
- Weight decay (weight_decay)
- warmup -epochs
- Warmup momentum
- Warmup initial bias
- Box loss gain (box)
- Class loss gain (cls)
- Class BCEloss positive weight (obj-pw)
- Object loss gain (obj)
- Object BCEloss positive weight (cls_pw)
- IoU training theshhold (IoU_t)
- Anchor multiple threshold (anchor_t)
- Focal loss gamma (f1-gamma)
- Images HSV-Hue , Saturation and Value augmentations (hsv_h, hsv_s, hsv_v)

- Image translation (translate)
- Image scale (scale)
- Image shear(shear)
- Image perspective
- Image flip up-down (flipud)
- Image mosaic(mosaic)
- Image mixup(mixup)
- Segment copy-paste (copy_paste)

There are many methods of setting these hyperparameters like grid search and random guessing. Here we will talk about a different approach to hyperparameter tuning, which is hyperparameter evolution.

3.11 Hyperparameter Evolution

Hyperparameter evolution is a Genetic algorithm-based approach to hyperparameter tuning. Genetic algorithm is a machine learning algorithm based on the concepts of natural selection. It is basically a probability-based search function that searches the optimal solution by mixing and matching parameters through multiple generations. The search begins by generating an initial population based on some given parameter. The chromosome or gene pool of the algorithm contain numbers that are a possible solution, new generations have created either crossover through mating or changed through mutation in order to go from one generation to another. Fitness is a function that helps define the optimal solution the algorithm is searching for. Fitness is evaluated on all members of each generation and only the ones with higher fitness are taken as the basis for creating the next generation via crossover or random mutation.[90]

In case of Hyperparameter tuning, the chromosomes are formed with hyperparameters of a model as genes.

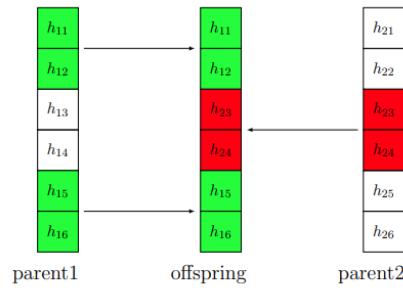


FIGURE 3.24: Hyperparameter evolution with GA

Some fitness function is determined based on the metrics of the task the model is meant to perform. Fitness is tested and new generations are created based on the GA. New generations are created using crossover of the highest-scoring chromosomes of the last generation. The mutation rate of the GA shouldn't be high as it may cause it to behave like an unintelligent random search function [91]. Based on the value of fitness scores the algorithm ends up with optimized values for all hyperparameters of the model.

CHAPTER 4

METHODOLOGY

In this chapter, we discuss how we implemented our proposed system. All implementations were done using Python 3.9.6 and its various modules and frameworks. Here we discuss our step by step procedures we followed in implementing the project.

4.1 Dataset Description

We took a more data-driven approach towards the classification of supermarket products. The dataset was created in different phases, and we will be comparing the models trained on different versions of the dataset later on in the report. Here we will discuss in detail the various versions of the datasets we built.

Primarily, we listed 16 of the most commonly found supermarket items. All of which are packaged products. The list of products we chose for our dataset are as follows:

1. ACI Pure Salt 1kg
2. Parachute Coconut Oil 50ml
3. Energy Biscuit
4. Dettol Original Soap

5. Mr. Twist
6. Teer Sugar 1kg
7. Rupchanda Cooking Oil 1L
8. Pran Hot Sauce
9. Cocola Chicken Masala Noodles
10. Frutika Grape 250ml
11. Sepnil HandSanitizer 40ml
12. Snickers
13. Sprite 250ml Bottle
14. Vaseline Original
15. Coca-cola 250ml bottle
16. Lux Soft Glow



FIGURE 4.1: Products in the dataset

We created the first version of the dataset was created with mostly single images of the products. Then we added synthetic images to enhance our dataset even further and finally, we added some stacked images to the datasets in order to improve detection performance.

TABLE 4.1: Versions of datasets

Dataset Version	Description
Version-1	Mostly single product images
Version-2	Synthetic Data was added to the dataset
Version-3	Synthetic Data was removed from the validation set
Version-4	55 Images were added with highly dense product arrangement

Throughout this chapter when we mention our dataset, we shall be mostly referring to Version-4. We shall discuss further the rest of the datasets in later chapters.

4.1.1 Data Acquisition

4.1.1.1 Daraz Review Image Scrapping

A comprehensive dataset consisting of our local retail products is rare to find. So we created our own dataset of retail products. In our dataset, we mostly accumulated data from 2 sources. Firstly, we have found that the most diverse images for a single product can be found in the review section of Daraz [92]. Daraz is an e-commerce site where people buy local goods. The review section of each product contains multiple images of the products they receive upon delivery.

We planned to capitalize on these images by creating a web-scraper that will iterate through each product review, download the review images, and collect them in a single folder. We made the web scraper in Python 3.9 using the modules Selenium 3.141.0 [93] and Beautiful Soup 4.9.3 [94].

Beautiful Soup is an HTML Parser. It has tools to retrieve the HTML data of a webpage. It requests the data from the webserver and stores them. The data is referred to as soup. However, there is a limitation that it can only retrieve the information a webpage displays on load. If some elements are loaded later via JavaScript may not be represented in the HTML data.

In order to get the review images to load before parsing HTML data, we needed to scroll down to the end of the reviews and keep moving to the following review pages. We achieved this by using Selenium which is a web driver. It drives a browser of choice as a human user would. It is possible to click buttons, scroll and perform all the interactions with a website like a human user.

We have used Selenium to drive Google Chrome on to load the product page and scroll down to the end of the current review section. After all the review images are loaded, we use beautiful soup to parse the HTML data and extract links to the images from that data. After that, we iteratively download the images from

those links. Then we scroll up and press the next button in the review section and scroll down again to load all the photos and keep following the same procedure until we can't find any more product images.

4.1.1.2 Manual Photography

We couldn't find all the products or enough images of some products from Daraz. So we manually photographed some of the products. The photos were taken using a smartphone camera against various backgrounds and lighting conditions. They were photographed both as individuals and as groups stacked together. We kept partial overlapping of products in some images but avoided any significant overlap for our model to properly recognize the product labels.

4.1.2 Data Cleaning

After acquiring the data, we had to clean the data, removing duplicates and unusable images. We had to deal with the fact that people didn't always upload exact pictures of the product in the review sections and some photos just focused on the expiry date rather than the entire product label. There were a lot of duplicate images as well.

We automated our duplicate image detection by using dupeGuru. It is a Python-based cross-platform tool for finding duplicate files. We provided a link to each product folder and the application automatically detected the duplicate files and removed them [95].

We have manually searched through the downloaded and photographed images for blurred images, images with products not correctly visible, images with significant overlapping, etc., and removed them from our dataset.

4.1.3 Pre-Processing

For us to train our model on the data, we had to first process the images. YOLO works great on square images. There is also a trade-off between training speed and

features when choosing the resolution of the images. If we choose a large resolution, it will be easier for the model to pick up more features while the training time will increase significantly. On the other hand, choosing too small a resolution will compress the features and lead to poor classification. So we decided to use 416x416 resolution for our model. But most of our images were not square. So we had to resize them. But if we wanted to crop them, we had to manually go through each image as the products were located in various parts of the images, and blindly cropping them could lead to data loss or the products being undetectable.

Thus, we resized the images by making them square first by padding the extra parts with 0s and creating a black region at the sides of the images. This was done using PIL, Matplotlib, and OpenCV. We wrote a script to transform all the images into 416x416 resolution JPG format.

4.1.4 Data Annotation

While there are many different ways for data annotation, we chose to annotate our data manually as in this way we can provide accurately labeled data to our model. So, we have used an annotation tool called LabelImg [96]. It's a graphical interface made using Python and PyQt4 for annotating image data in different formats. It also supports the YOLO format.

We simply had to put our list of classes into its data folder and using its interface, we have drawn bounding boxes around products in each image and selected their class. Then it automatically generates the text file in the format mentioned above.



FIGURE 4.2: Correct vs Incorrect Annotation

It's worth mentioning that we have tried our best to make an exact bounding box around the object, not including extra areas as it can hamper the model's accuracy. One of the most crucial accuracy metrics for an object detection model is IoU or Intersection over the union. If we have extra spaces in our bounding box, this metric decreases.

4.1.5 Synthetic Data

We chose to augment our dataset with synthetic data in order to enhance the overall performance of our model. In this section we have discussed the techniques for synthetic data generation we have used.

In order to create synthetic images, we extracted transparent PNG images of each object (both front and backside) using Adobe Photoshop. Then we collected around 20 different surface textures and placed products randomly on each of the background images. The position, rotation, and scaling of the products were randomized. Each synthetic image had 1-10 products in them which in turn lead the number of synthetic instances to far surpass the number of authentic instances in the dataset.

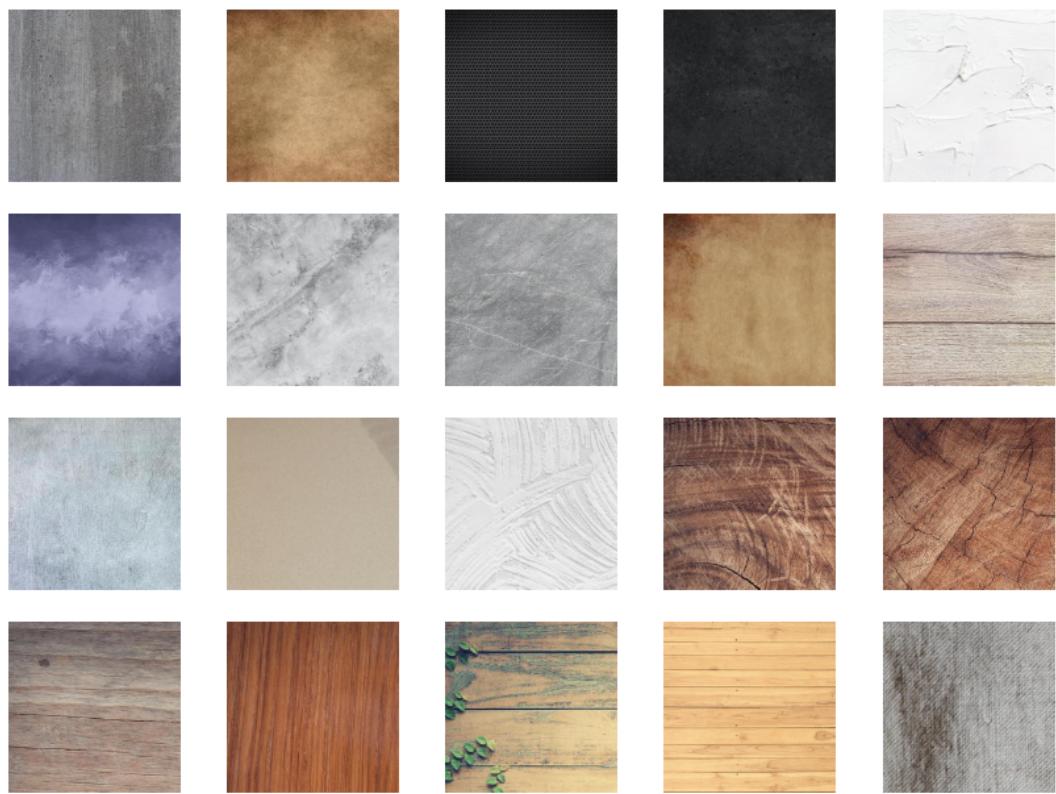


FIGURE 4.3: Surface textures for synthetic image generation

It requires the background and object files to be in separate folders. Each object category needs to be in its subfolders. The script then randomly chooses some number of objects and loads the object PNG files and randomly translates, rotates, and scales them across a randomly selected background image. It creates the specified number of synthetic images in the same way.



FIGURE 4.4: Synthetic Image

This placement of images was done using the Flip Github repository. It takes the PNG images of the objects and the background images as input and generates any specified number of random Synthetic images. There is some discrepancy in the scale of the images relative to each other as no relative size ratio is fixed provided in the algorithm. But this doesn't hurt the accuracy of the model very much, rather the inclusion of different sizes of the products helps it to be able to better recognize the product from different distances. We initially have created 150 synthetic images per class.

4.1.5.1 Annotating Synthetic Data

The annotations for the bounding boxes in the Synthetic images were automatically created on generation. Flip used Pascal VOC format for annotation so we have to convert the format to YOLO text files before using them [97]

Pascal VOC format uses XML files to store annotation details. There are extra parameters in this format than what YOLO needs like name, truncated, pose, difficult, etc also the bounding box format for Pascal VOC is different (x_{min} , y_{min} , x_{max} , y_{max}). So, in order to convert the annotation to YOLO format, we must convert the bounding box to Yolo format and replace the object name with its id number. [98]

$$\begin{aligned} norm_x_center &= \frac{(x_max + x_min)}{2 * img_width} \\ norm_y_center &= \frac{(y_max + y_min)}{2 * img_height} \\ norm_width &= \frac{(x_max - x_min)}{img_width} \\ norm_height &= \frac{(x_max - x_min)}{img_height} \end{aligned}$$

We created a python script for automatically converting the Pascal VOC XML files to YOLO Format text files. After running the conversion the synthetic images have been ready to be added to the dataset

4.1.6 Dataset Distribution

The dataset was initially split between with an 80, 10, 10 ratio between the train, validation, and test slices. Then the synthetic data was added in the aforementioned process. We have 7348 object instances spanning across 3430 images. The instance distribution among the 3 slices for the final version of the dataset are shown on Figure 4.5. And all synthetic values for the training dataset have been placed only in the training slice. We have a total of 1898 synthetic images in the dataset. The distribution of authentic and synthetic data in the training slice is illustrated on Figure 4.6.

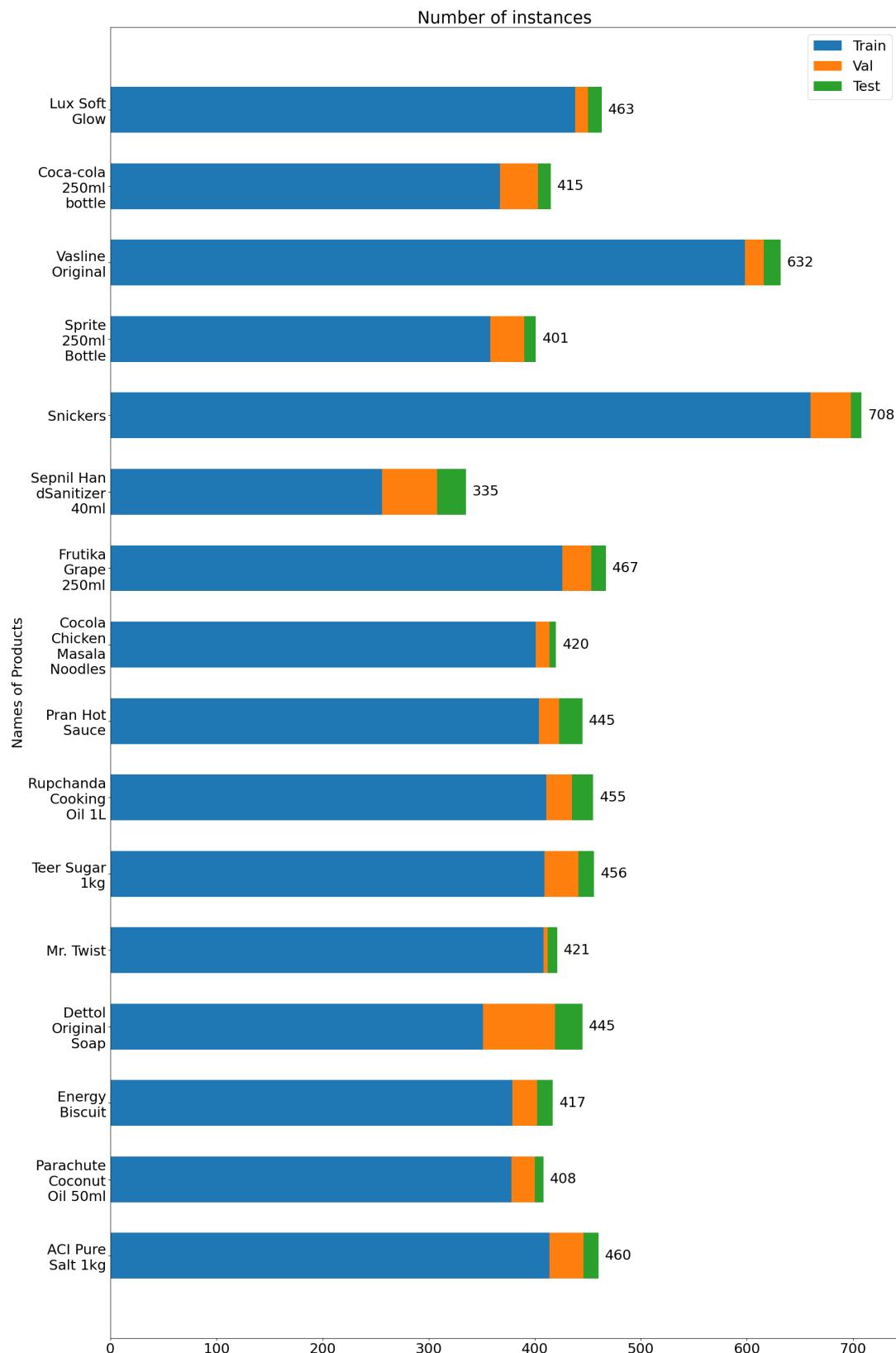


FIGURE 4.5: Distribution of products in dataset splits

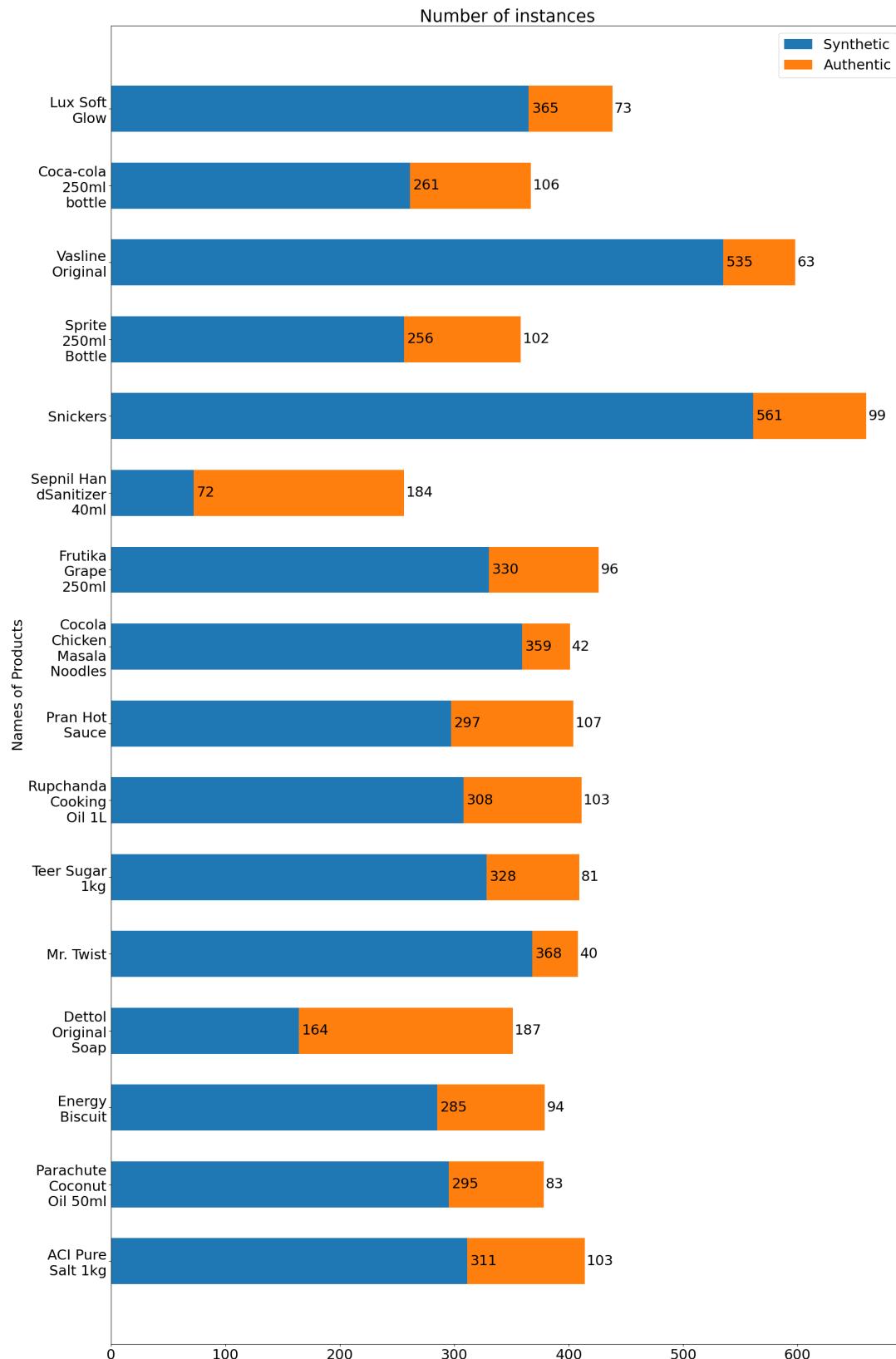


FIGURE 4.6: Distribution of Synthetic and Authentic images in training set

4.2 Model Overview

The algorithm YOLOv5 has been implemented with Pytorch. It's a machine learning platform that emphasizes both the aspect of speed and usability. Pytorch is based on Python, so it is compatible with all the renowned libraries for scientific computation. It also supports GPU acceleration for maximum performance for maximum computing power. [99]

Yolov5 has multiple variations of models depending on size and complexity. There are 4 variations:

1. Yolov5s
2. Yolov5m
3. Yolov5l
4. Yolov5x

We took a transfer learning approach by loading pre-trained weights that were trained up to 300 epochs on the COCO dataset. Then trained the models on our dataset starting from the pre-trained weights as a checkpoint and trained up to 100 epochs. Then saved the best weights that we obtained from the training.

COCO or Common Objects in Context is a dataset by Microsoft which features a wide range of common object instances. The dataset is designed for object recognition with an understanding of the entire scene. There are 2.5 million instances of common everyday objects spanning across 328,000 images that are labelled for bounding boxes and object instance segmentation[100].

The pretraining information on the COCO dataset for the YOLO models are showcased below [101]:

TABLE 4.2: Pretraining information on the COCO dataset

Model	Size (pxl)	mAPval 0.5:0.95	mAPtest 0.5:0.95	mAPval 0.5	SpeedV100 (ms)	FLOPs 640(B)
s	640	36.7	36.7	55.4	2.0	17.0
m	640	44.5	44.5	63.1	2.7	51.3
l	640	48.2	48.2	66.9	3.8	115.4
x	640	50.4	50.4	68.8	6.1	218.8

We only trained the Yolov5s and Yolovm models due to GPU resource limitations. We can infer from the graph below that larger models need more GPU processing power. We figured that the optimum point for our available hardware would be Yolov5m or medium model.

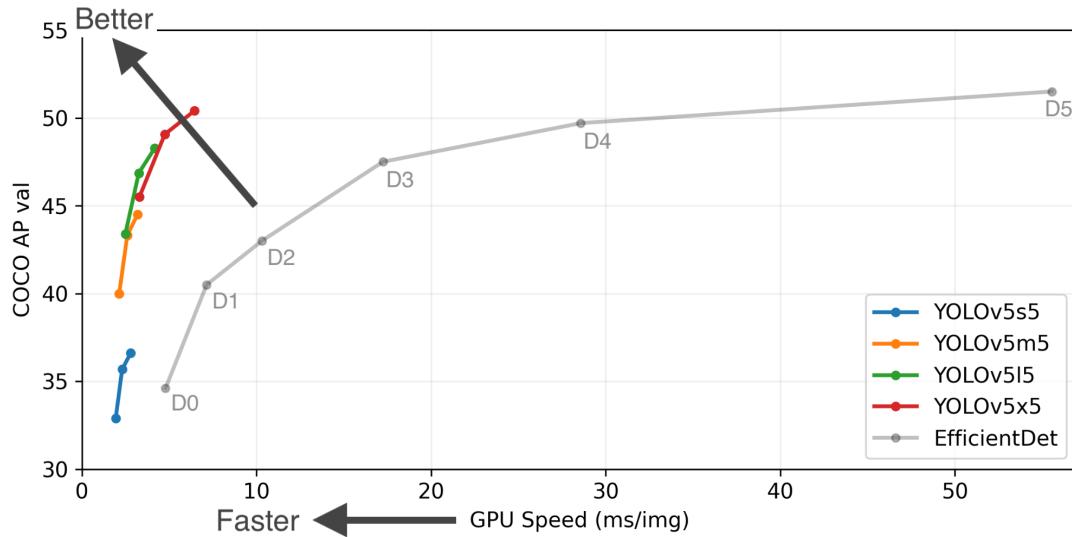


FIGURE 4.7: Comparison of Yolov5 Models

The difference between the small and medium models is that the small one has layer channel multiple and model depth values of 0.5 and 0.33 compared to 0.75 and 0.67 of the medium model. They have similar overall structures. The structures for the two models are show on Figure 4.8



FIGURE 4.8: Yolov5 Architecture [10]

4.2.1 Hyperparameter Evolution

We have used a Genetic algorithm-based hyper-parameter evolution technique to tune our model further. Each generation to 10 epochs starting from the default hyperparameter, which is optimized on the COCO dataset.

The fitness function is defined as:

$$F = 0.1 * (mAp_{0.5}) + 0.9 * (mAp_{0.5-0.95})$$

This algorithm has an 80% probability of mutation with a 20% variance. Instances are created until the mutations occur. Then the best mutation of each generation is used to create the next generation in the same way. At last, the hyperparameters for the best generation are saved.

We trained a total of 150 generations and ended up with the best hyperparameter settings on the 11th generation. Even though it was recommended to run 300 generations, we stopped early due to lack of GPU power.

4.2.2 Training Description

It took around 40s per epoch ending in around 1hr of total training time. The training description for Yolov5m model is listed below:

Hardware:

TABLE 4.3: Hardware used to train the models

CPU	Intel i7-8700
RAM	32 GB
GPU	MSI GTX 3080 Trio

Settings:

TABLE 4.4: Training settings

Epoch	100
Batch Size	32
Image Size	416x416
Workers	1
Optimizer	SGD

4.3 Hardware Setup

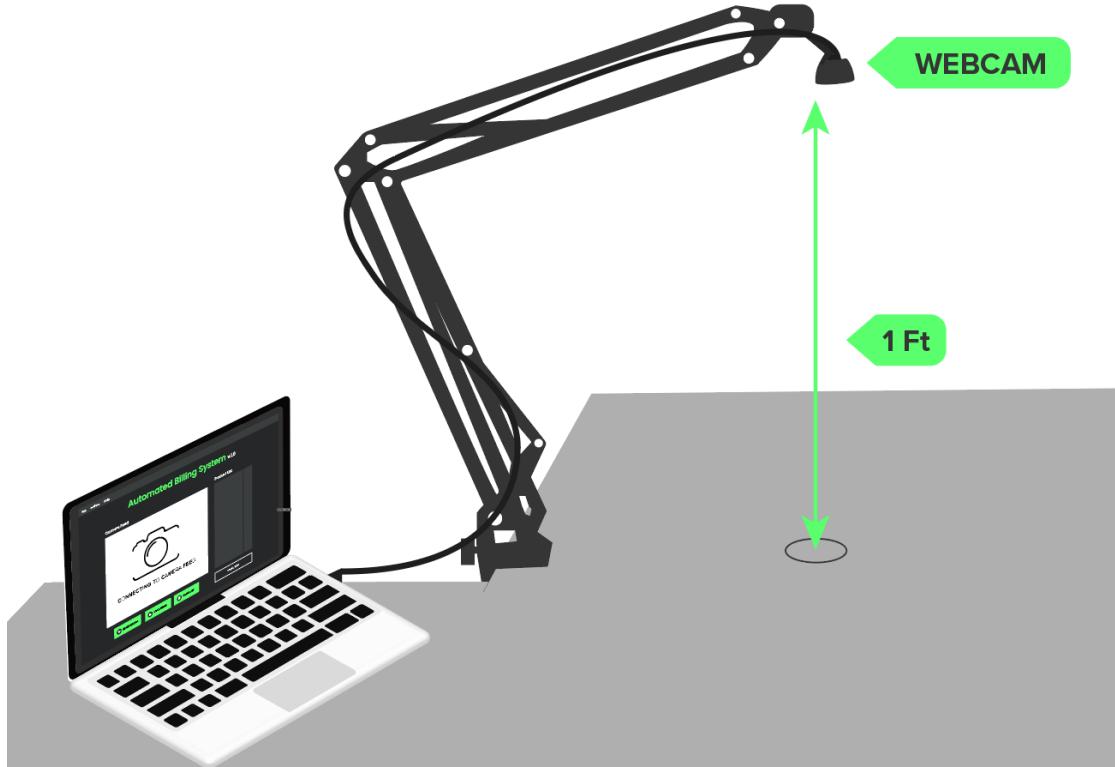


FIGURE 4.9: Comparison of Yolov5 Models

We constructed a table-mounted setup with our components. The setup consists of the following:

1. Webcam
2. Adjustable Table mounted arm
3. Computer
4. Table

The webcam is kept approximately 1 foot above the surface of the Table using the adjustable table-mounted arm. The surface of the Table should have plain colours devoid of complex patterns to aid with the detection. The webcam is connected to the computer via a USB cable, as shown in the figure above. Products will be placed in batches under the webcam for detection and billing.

4.4 GUI Implementation

GUI or Graphical User Interface is an interaction mechanism where a user can use a pointing device (perhaps a mouse) to interact with different functionalities of a program. This brought about a revolutionary change in the field of Human-Computer Interaction (HCI) in modern times. We planned to implement a GUI base to provide the cashier with an easy-to-use application for our billing system [102].

4.4.1 Objectives of the Application

The list of features we have considered implementing are as follows:

1. Initiating and terminating a billing session
2. Clearing existing list
3. Generating a printable receipt
4. Dynamically updating prices to a database
5. Automatic refreshing of the product list when products are changed
6. Mechanism to lock an existing batch of products on the list

Based on the features, we have first created a wireframe of the GUI layout.

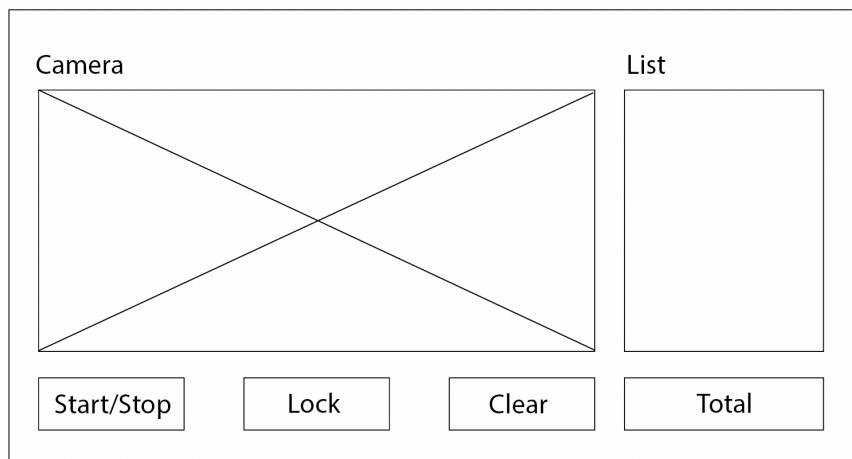


FIGURE 4.10: Wireframe of the app

4.4.2 Description of the UI

The UI design was done based on the wireframe diagram, using Adobe Illustrator. Then the GUI was implemented using the PyQt5 framework [103].

Qt is a cross-platform C++-based programming toolkit that offers everything to build an app UI. PyQt is the collection of bindings that glue together Qt and the vast array of modules that we get access to with Python.

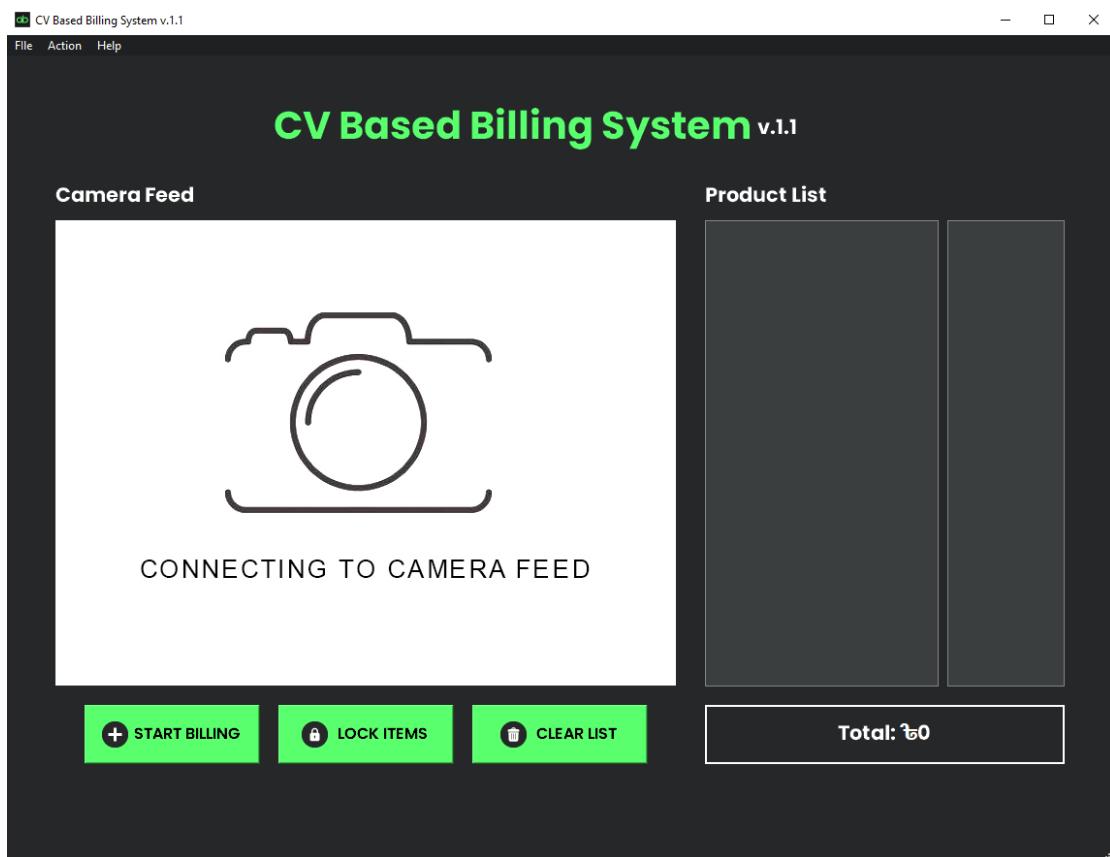


FIGURE 4.11: GUI of the Application

Camera Feed: In the app, the Camera Feed window displays the camera view in real-time. Any products that are detected are marked with a bounding box around them, and their name is displayed.

Product List: When the billing starts, the detected products are listed here. The first column shows the name and amount of the product in this format: n x Product_name. And the second column lists the value of those products in Taka.

Buttons:

- The “Start Billing” button initiates the billing process. The application won’t list the products before the button is pressed. It’s a toggle button, so it changes to “Stop billing” after billing starts. To stop the billing process, we have to press the same button again.
- The “Lock Items” button is for locking down the batch of products currently on the list.
- The “Clear List” button clears all the current entries in the list.
- The “Stop Billing” button terminates the billing process and generates a receipt in a text file.

All of the functionalities are tied to keyboard shortcuts for ease of use of the cashier.

TABLE 4.5: Shortcuts for the App

Functionality	Keyboard Shortcut
Start billing	Space
Lock items	S
Clear list	X
Stop billing	Space

4.4.3 Highlight Features of the App

4.4.3.1 Automatic Refresh

Usually, no object detection system is 100% accurate under every scenario. Thus, the application may momentarily miss-classify the products when the cashier places them on the checkout area, or one product is significantly overlapped by the other.

The application automatically updates the list on each frame to tackle this issue. So that if a product is miss-classified on one frame, the application will dump the

data and refresh it on the next frame automatically, or the cashier can slightly rearrange the products on demand to aid it classifies the products correctly.

4.4.3.2 Batch Locking Mechanism

One of our main objectives is to make the process of supermarket checkout faster by enabling the cashier a way to bill multiple products at once rather than one at a time. If we could bill all the products together, then it would have been most efficient. But there is no practical way to achieve that as cameras have a limited Field of View (FOV). We circumvent this limitation by introducing a way to bill the products in batches.



FIGURE 4.12: CV based billing system

After the billing starts the detected products are automatically added to the list. As we automatically refresh, the products will disappear from the list if we

remove them from the checkout area. In order to achieve the batch billing system, we introduced the locking mechanism. When “Lock items” is pressed, it locks the current items to the top and highlights them in green. We can see in the Figure above that the items highlighted in green are locked. When we finish billing, only the locked items in the list will be billed.

Apart from the locked items, we can see the same items in a non-highlighted state under it. The application is considering the item as a new entry; if we remove the item from view, it will automatically be removed. Furthermore, if we press stop billing in this state, it will only bill the locked items.

4.4.3.3 Receipt Generation

A receipt is generated after “Stop Billing” is pressed. All the locked items are taken from the list and presented in a POS-style receipt format in a .txt file. This is just for the representation of the receipt. In a practical application, we plan to convert the .txt format to POS commands in order to print the receipt via a POS thermal printer. The prices for the individual items are stored in a local CSV file. One can change the values from the file in order to update the prices.

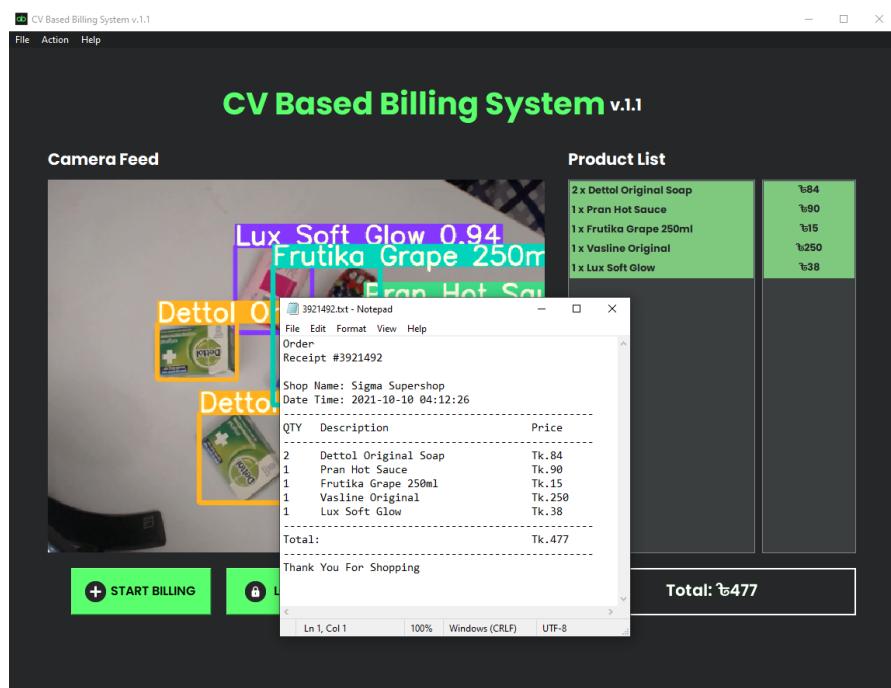


FIGURE 4.13: Generated Receipt

CHAPTER 5

RESULT AND ANALYSIS

In this chapter, we aim to discuss the results of our implementation and analyze its performance in various scenarios and test setups. We also discuss how we finalize the model to use as the backend of the billing system application.

5.1 Dataset Analysis

Let's analyze the bounding box data from the last version of the dataset. We use correlograms as a method of analysis. Seaborn pair plots [104] create correlograms by creating a histogram plot with each pair of parameters. The plots at the diagonal are simple distributions.

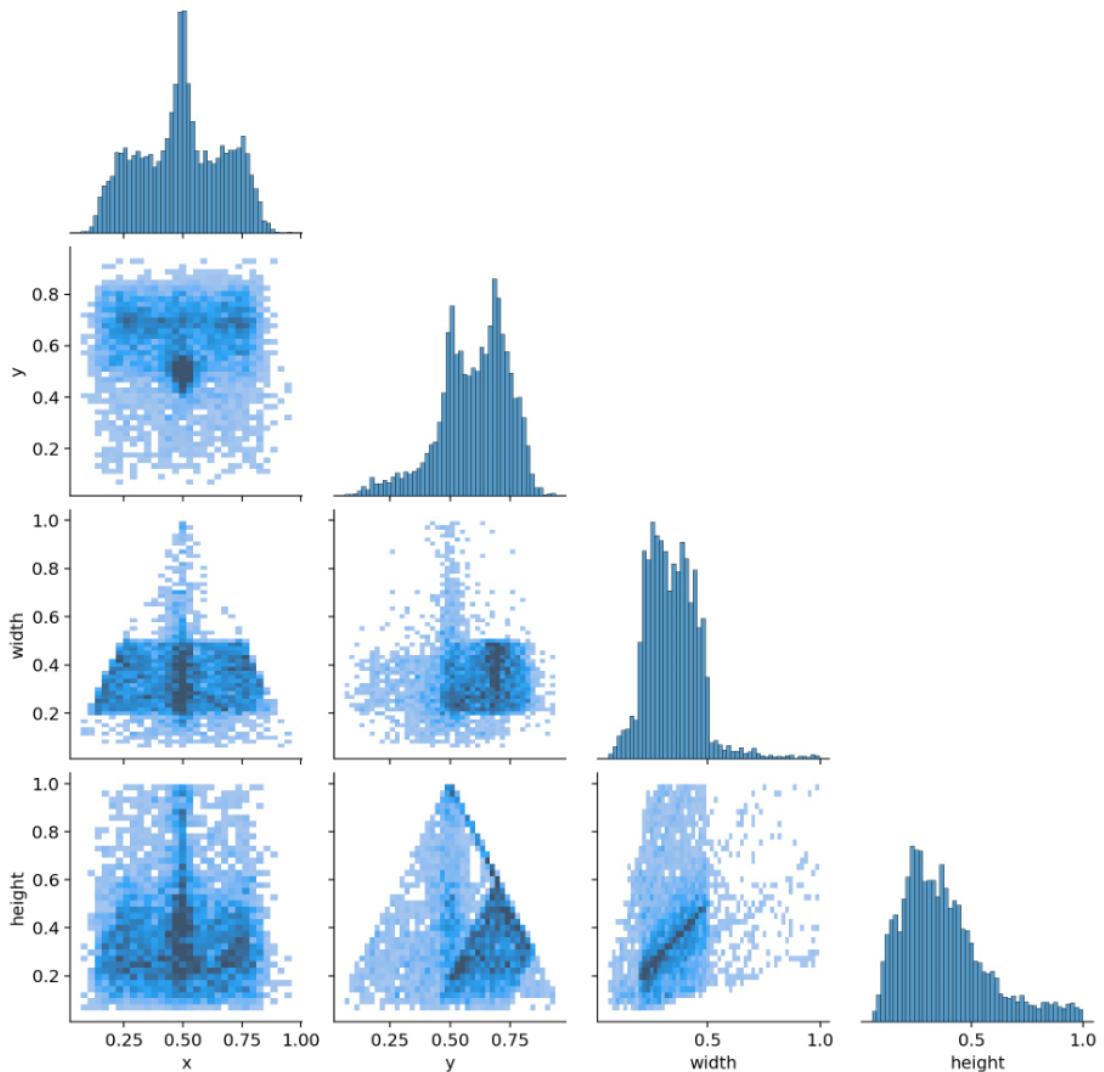


FIGURE 5.1: Correlogram of the dataset

These are pair plots of the bounding box parameters of each of the object instances. The only graphs of interest are xy and height-width graphs. It gives us an idea about the distribution of the objects throughout the images. From xy figure, we can infer that a lot of the images have x and y coordinates near the center of the image as there is a very dense spot in the middle of the plot. This is expected as there are a lot of daraz review images in our dataset which are just the picture of the single product. From the height-width plot, we can infer that there is a good amount of square bounding boxes in the images as there is a dense diagonal line that represents a linear increase in height with an increase in width.

5.2 Evaluation Metrics

We will discuss some popular metrics for object detection in this section. There are two parts of the object detection process, the first is object localization and the other is classification. One of the used evaluation metrics for localization is Intersection over Union (IoU)

5.2.1 Intersection Over Union (IoU)

IoU is the ratio of the area of intersection between the labeled bounding box (Bl) and the predicted bounding box (Bp) and the area of the union of them. A higher value of IoU represents a more correct localization of an object. For an application, usually, a threshold value is selected and if IOU is greater than that threshold it is considered to be correct [105].



FIGURE 5.2: Intersection over union

In the Figure 5.2, the green bounding box is the ground truth or labeled bounding box and the yellow one is predicted by the model. The shaded area represents the intersection between them. IoU is defined by the ratio of the shaded area and the sum of the areas of the two bounding boxes.

5.2.2 Evaluation of Object Classification

After localization comes the classification of objects. Before we begin some common concepts should be discussed, **True Positive**: If the object has been correctly predicted in a scenario where it is present. **True Negative**: If the object has been correctly not found in a scenario where it is absent. **False Positive**: If the object has been detected in a scenario where it is absent. **False Negative**: If the object has not been detected in a scenario where it is present.

Confusion Matrix:

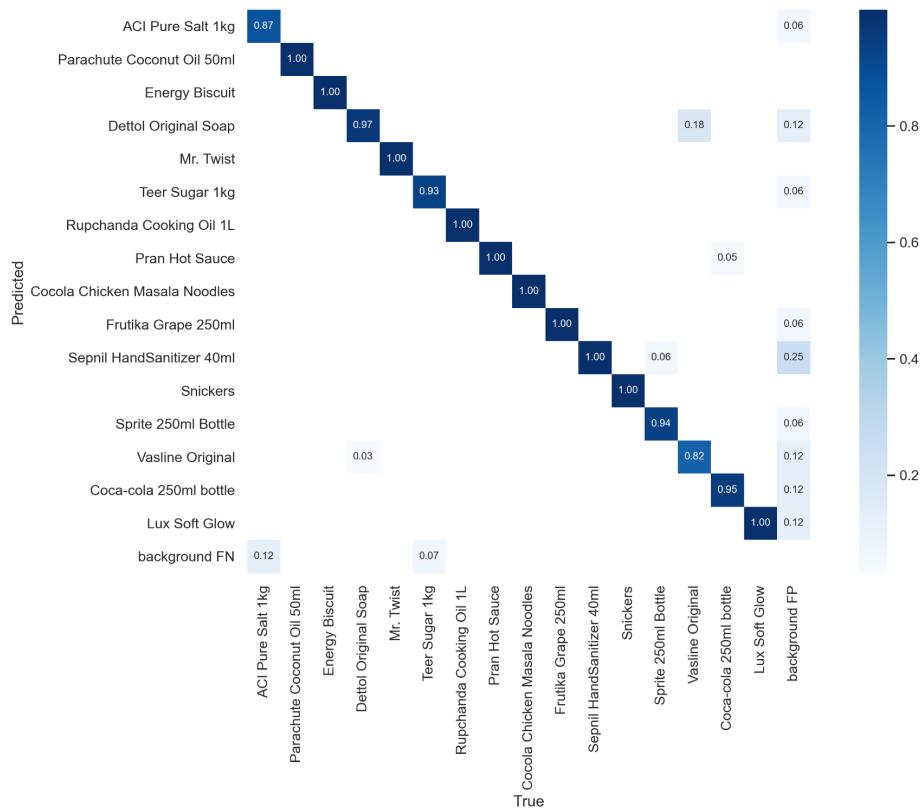


FIGURE 5.3: Confusion Matrix

A confusion matrix depicts the points where a model is confusing between classes. It is a nxn matrix with true and predicted results at perpendicular and lateral sides. The diagonal of the matrix represents the true positive predictions. Other points of the matrix show where the model has misclassified. For example in the figure above, 18% of Vaseline Original instances were confused with Dettol Soap. It is a very effective way of evaluating a model on a dataset [106].

Precision:

Precision is the measurement of how many positive predictions are actually correct. It focuses on the correctness of positive detection which can be a good metric for cases where a false positive can be more troublesome than a false negative.

$$\text{Precision} = \frac{TP}{FP + TP}$$

Recall: Recall is the measurement of how many positive predictions were done in scenarios where the object was actually present. It is a good metric where false negatives are more troublesome.

$$\text{Recall} = \frac{TP}{FP + FN}$$

F1 Score: It is the combination of Precision and Recall scores. It is a good choice for tasks that are sensitive to both false positives and false negatives.

$$F1Score = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

P-R Curve & Average Precision(AP): Precision-Recall Curve can be seen as a balance between the two metrics. For a good model, the precision should increase with the increase in recall. The area under the curve is a good metric for evaluation. One of the techniques for calculating the area under the curve is done by 11 point interpolation method. The precision value for 11 uniformly distant recall levels are averaged to get the value for AP [107].

$$\text{AP}_{11} = \frac{1}{11} \sum_{R \in \{0.0, 0.1, \dots, 0.9, 1\}} P_{\text{interp}}(R) \quad (5.1)$$

where,

$$P_{\text{interp}}(R) = \max_{\tilde{R}: \tilde{R} \geq R} P(\tilde{R}) \quad (5.2)$$

$P(R)$ where the value of recall is greater than R is used for AP calculation instead of precision at R value.

Mean Average Precision(mAP): mAP is the mean of the AP for all classes in a dataset, it is one of the most used metrics for object detection evaluation.

$$\text{mAP} = \frac{1}{N} \sum_{i=1}^N \text{AP}_i \quad (5.3)$$

Most of the renowned competitions use this metric as their main metric such as the PASCAL VOC Challenge, COCO Object detection challenge [108] and the Open Images challenge [109].

Confidence Score: It is the probability of the predicted object being in the bounding box.

Accuracy: It is defined by the total correct predictions by total predictions. It is not an ideal metric for object detection as the dataset may not be balanced

5.3 Test setups

We have created a total of 3 test setups each focusing on different aspects of the object detection problem.

Test Set 1: This test set focuses on the highly dense placement of the products and how our models perform on them. We are using 30 images containing a total of 351 object instances. This set of images are arranged in two parts, front sided images and back-sided images. This is used for evaluation under 2 different lighting conditions. One where it is well lit and another where there is not enough light.

The dark setting was created artificially by lowering the brightness of the images in order to keep other factors unchanged.

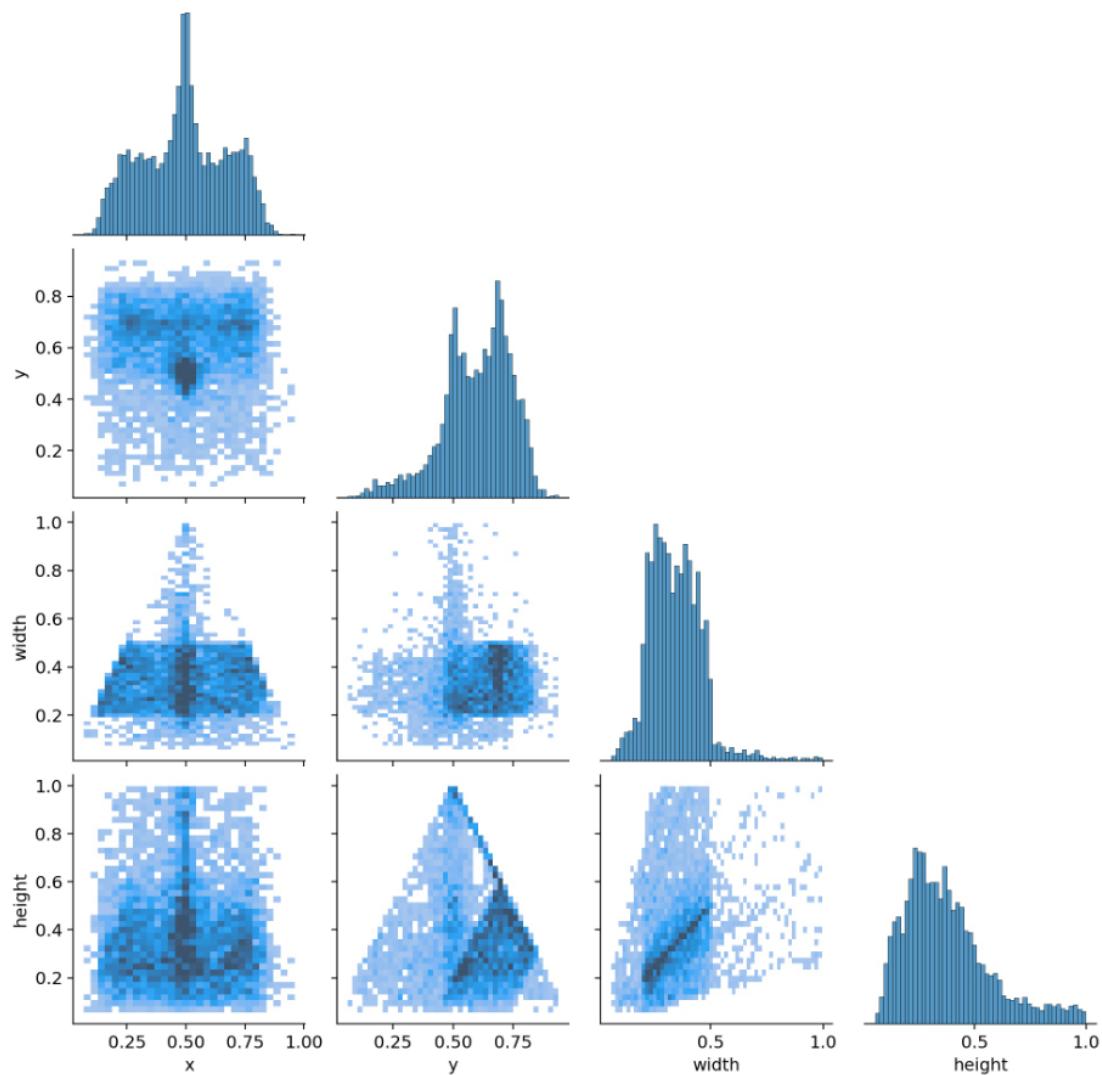


FIGURE 5.4: Samples from Test Set - 1

The distribution of front-facing products in the set are shown in the pie chart below:

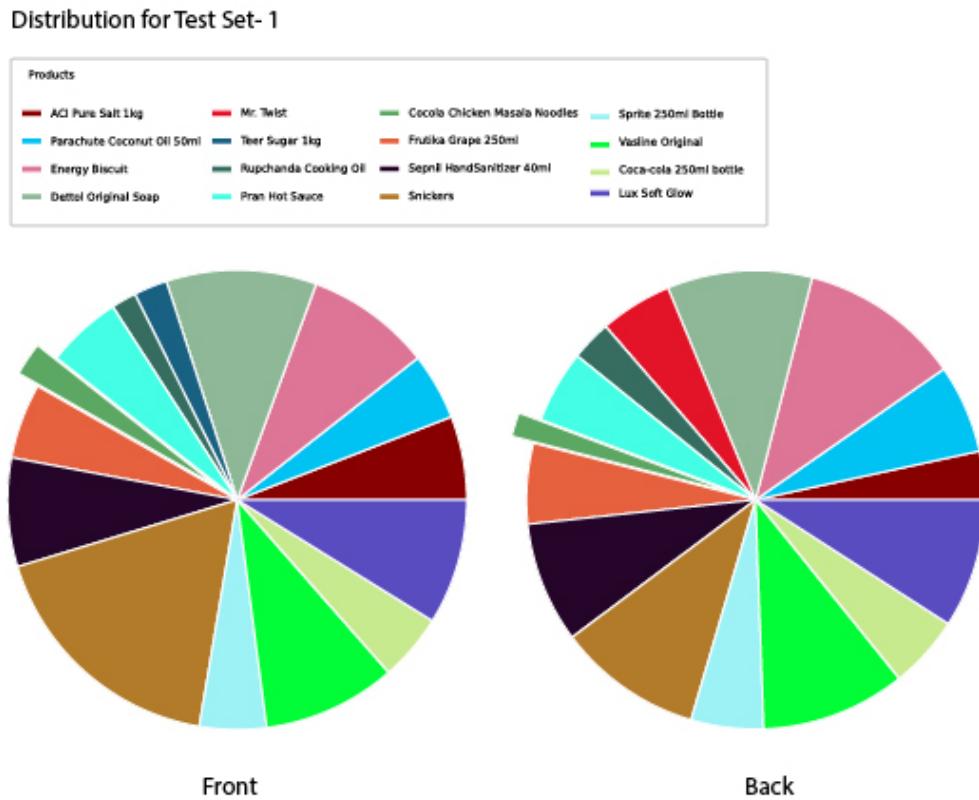


FIGURE 5.5: Distribution of Test Set -1

Test Set 2: This test set focuses mainly on real-time detection. This is a 1:30 min long clip of a person using the system to bill products. We take several frames from different locations of the clip to evaluate our models.

Test Set 3: This test set focuses mainly on the detection of products on different backgrounds; it contains a 10% split of the original dataset. The products on this set are on various backgrounds allowing for a robust background-independent test.

The distribution of front-facing products in the set are shown in the pie chart below:

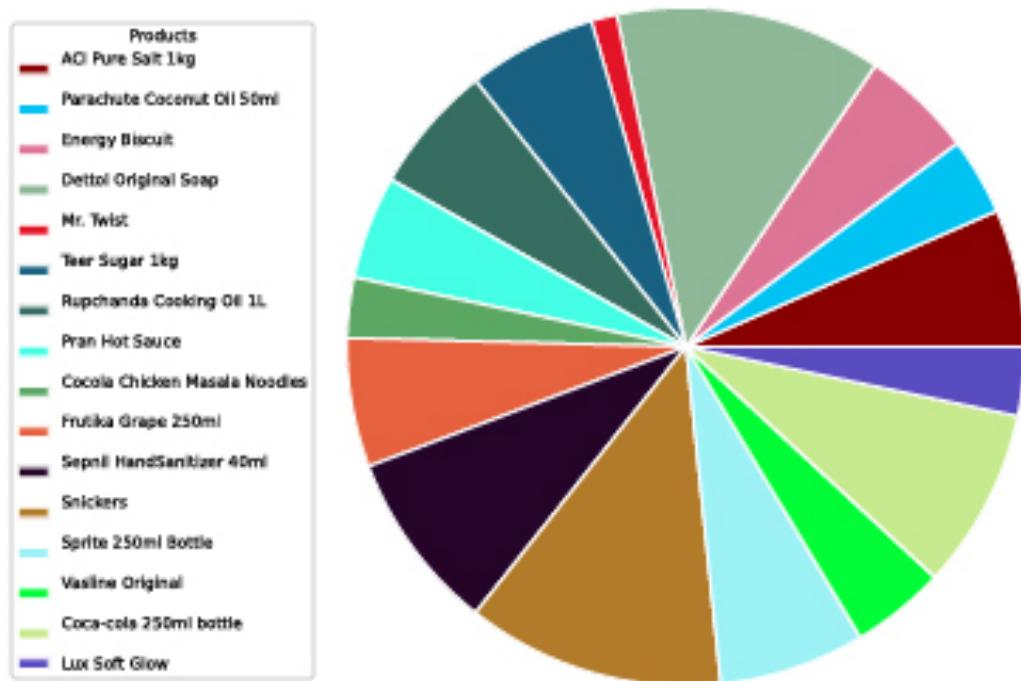
Distribution for Test Set- 3

FIGURE 5.6: Distribution of Test Set -3

5.4 Models for Evaluation

In different stages of development, we have trained different models either based on different model architectures or different states of the dataset. So far, we have only talked about our final result. In this section, we shall specify all of the models we used for testing:

1. **Synthval_Small:** Yolov5s model. Our validation set contained synthetic images during this phase, and most of the dataset was full of individual product images.
2. **NoStack_Small:** Yolov5s model. We removed synthetic images from our validation set.

3. **Stacked_Small:** Yolov5s model. We added more densely packed multi-product images to our dataset
4. **Stacked_Medium:** Yolov5m model. We added more densely packed multi-product images to our dataset

5.5 Training Results

Here we showcase the training statistics of the stacked medium model. There are a couple of parameters here.

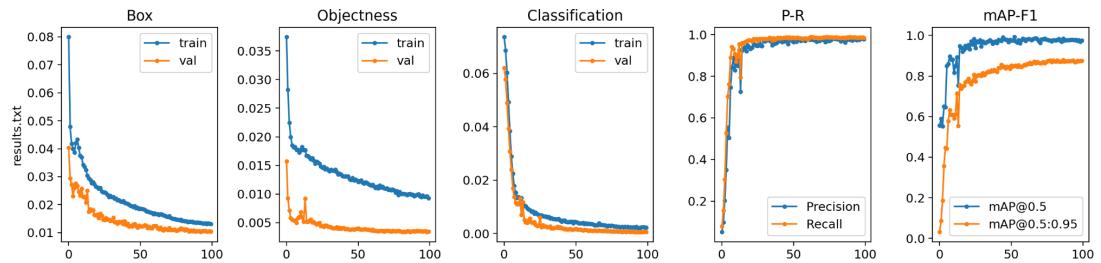


FIGURE 5.7: Training Statistics for Stacked_Medium model

The Box metric is actually the coordinate loss; it helps the model predict better boxes. The objectness actually represents the objectness loss, minimizing it helps the model improve IoU. The Classification loss does a similar thing in the case of improving classification. These curves are plotted for both the training set and the validation set. We can notice that the losses are slowly decreasing as epochs increase. Precision and Recall curves are also given. We can notice that both of them increase with epoch numbers and reach a saturation level near the end. The mAP curve also follows a similar pattern, it starts off a bit zig-zag, but it eventually reaches a high value and saturates. The training statistics for most of the models look similar; that is why we are only showcasing our overall best model.

5.6 Results on Densely Arranged Objects

The models mentioned above are all evaluated on Test set - 1. We can illustrate the results for each model using confusion matrix below:

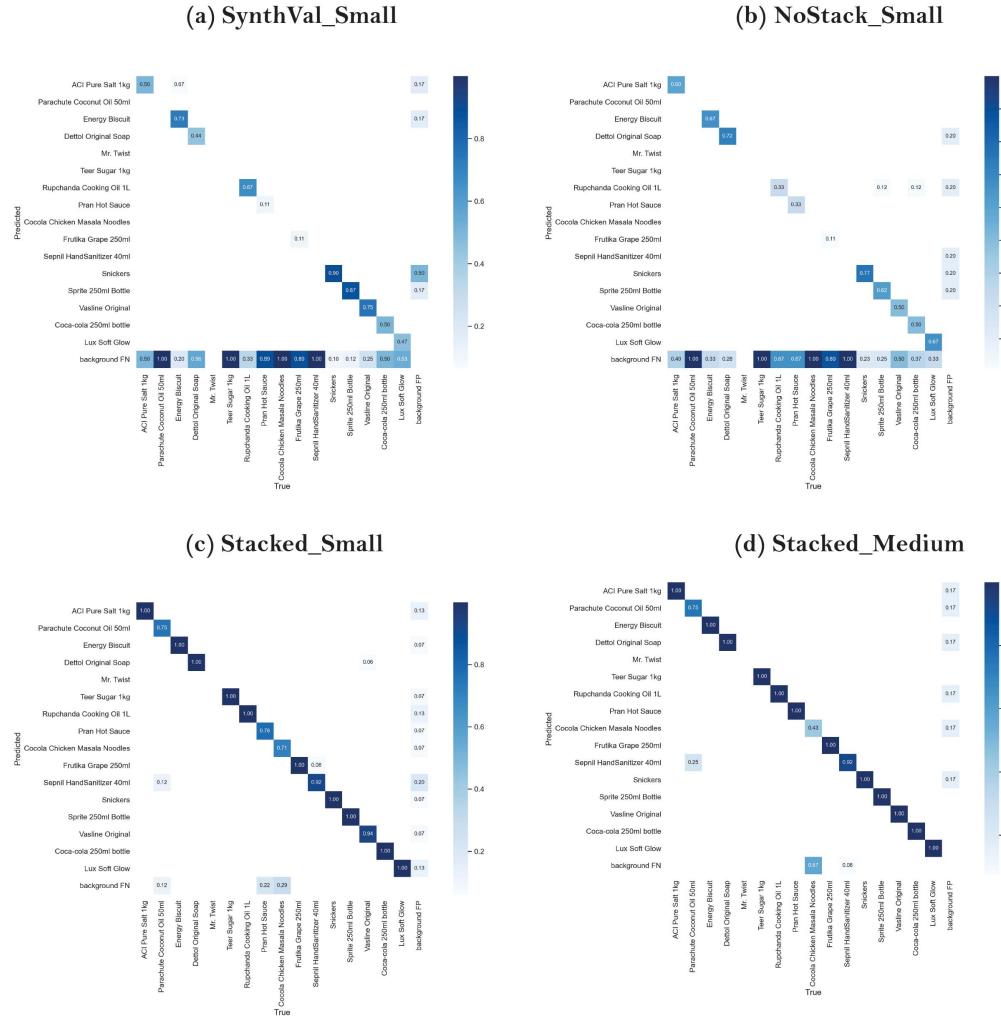


FIGURE 5.8: Confusion Matrices for Front facing products in Test set-1

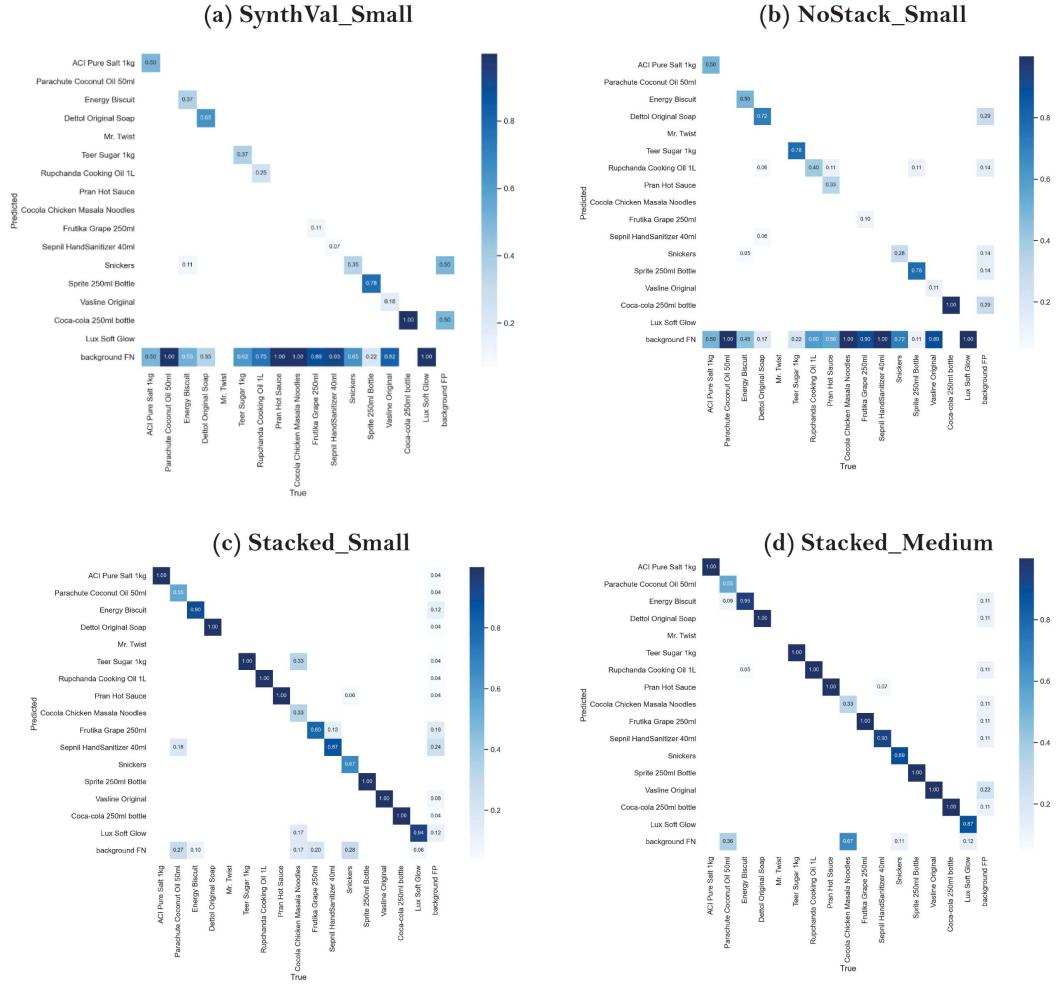


FIGURE 5.9: Confusion Matrices for Back facing products in Test set-1

If we notice the matrices (a) and (b) for both figures (front and back) we see a noticeable improvement. This is the result of removing all the synthetic data from our validation set. The first model was fitting more towards synthetic images as they were in the validation set. Thus reducing its accuracy on real images that are often not as perfect in terms of background separation as synthetic ones.

More images of the densely packed product were added between (b) and (c). Here we can see substantial improvement in detection It has gotten very good at

correctly classifying the products at this point.. Also between (c) and (d) can observe more improvement as we move to a larger model which is Yolov5.

There is one important thing to mention here is that from the Pie charts of the test set 1 we can see that there is no representation of Mr.Twist in this set, hence the discontinuity in the diagonal of the matrix. Also, the reason for Cocola Chicken Noodles to have the lowest score across the models is that the item has a very low presence in this set. We can also notice that the performance of the models is better on the front labels than on the back. The prime reason for this is the fact that data of front and back images of each product are not distributed uniformly in the dataset. The front-facing images are most likely to be greater.

5.7 Results in Different Lighting Conditions

mAP scores for the 4 models are shown in the chart below for both front and back labels, each of which is in bright and dark settings separately.

In the Figure 5.10, each group of 8 bars represent mAP scores for each model. The first subgroup group of 4 bars represent performance on front labels and the other 4 on back labels. First two bars for each subgroup represents mAP@0.5 for bright and dark lighting conditions, the other two represents mAP@0.5:0.95 scores for bright and dark lighting conditions within the given test set slice.

If we analyze the information in the chart, 3 points of improvement can be noticed. Firstly, we can see that the overall mAP scores improved a bit after the removal of synthetic data from the validation set. Secondly, the scores significantly improved after the inclusion of the 55 highly dense images of products in the dataset. Thirdly, we can notice a bit of improvement after changing to a larger model.

The scores for the back labels are somewhat lower than the front labels scores. There are two reasons that might have caused this. The back labels of most products don't have enough recognizable features, also the number of backside images are low in the dataset as most authentic images come from Daraz reviews.

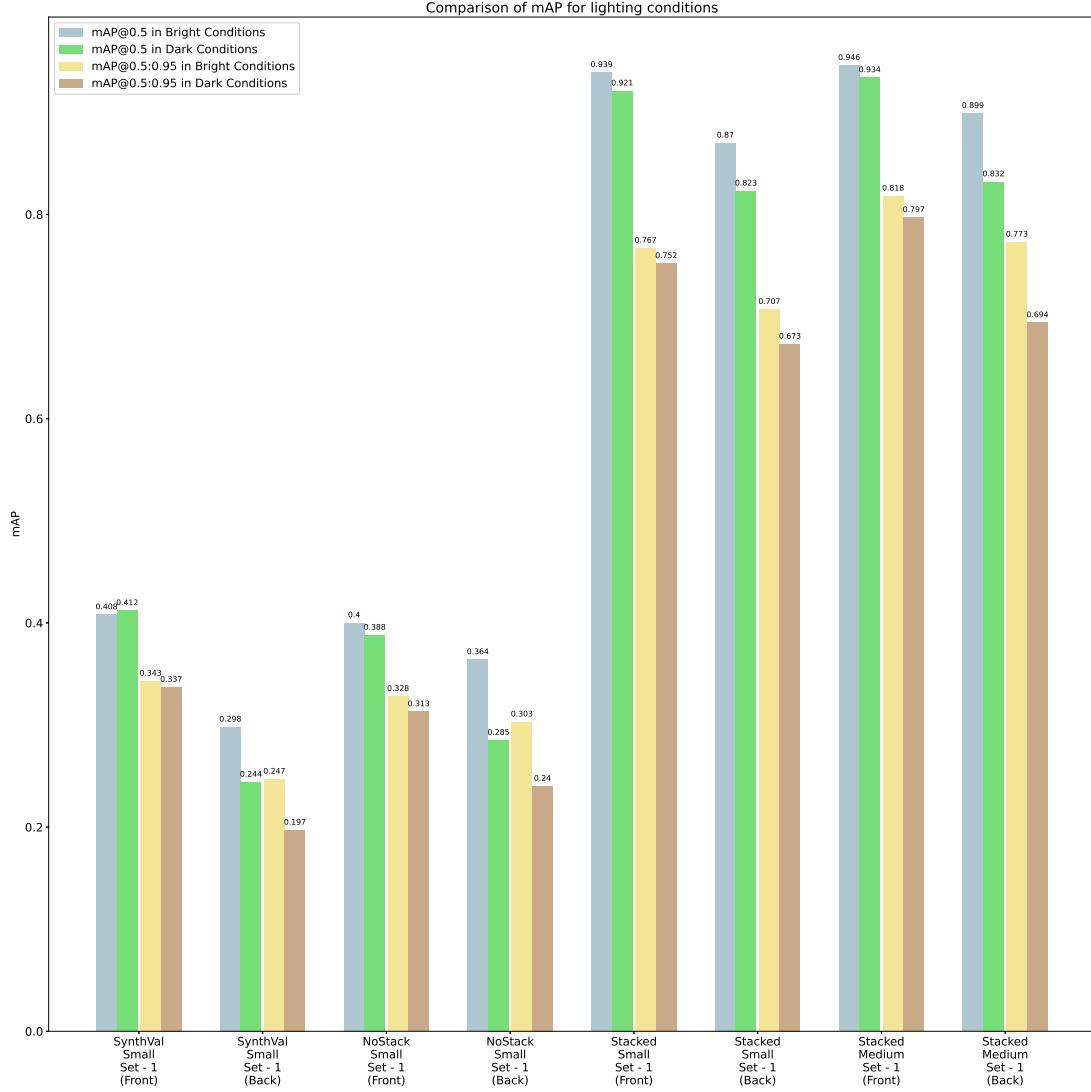


FIGURE 5.10: Comparison of mAP in different lighting conditions

5.8 Analysis of Real-time detection

Even though it's difficult to put the results of real-time detection into a representable form on paper. We make an attempt to evaluate performance based on our Test Set - 2. We pick several frames from the video and evaluate performance based on confidence scores and the number of false positives and false negatives. This depicts how confidently the model is detecting the products in each scene. The scenes are shown below



FIGURE 5.11: Selected Frames for the Test-2

The confidence scores and number of errors for each of these scenes are plotted in bar charts below. The evaluation of performance should depend on both of the parameters together.

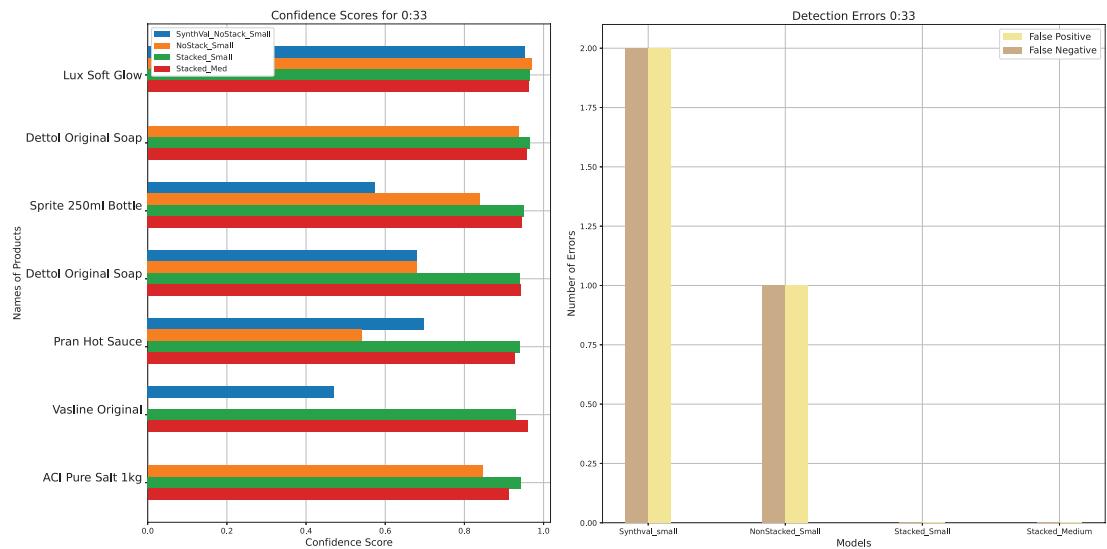


FIGURE 5.12: Analysis for time frame 0:33

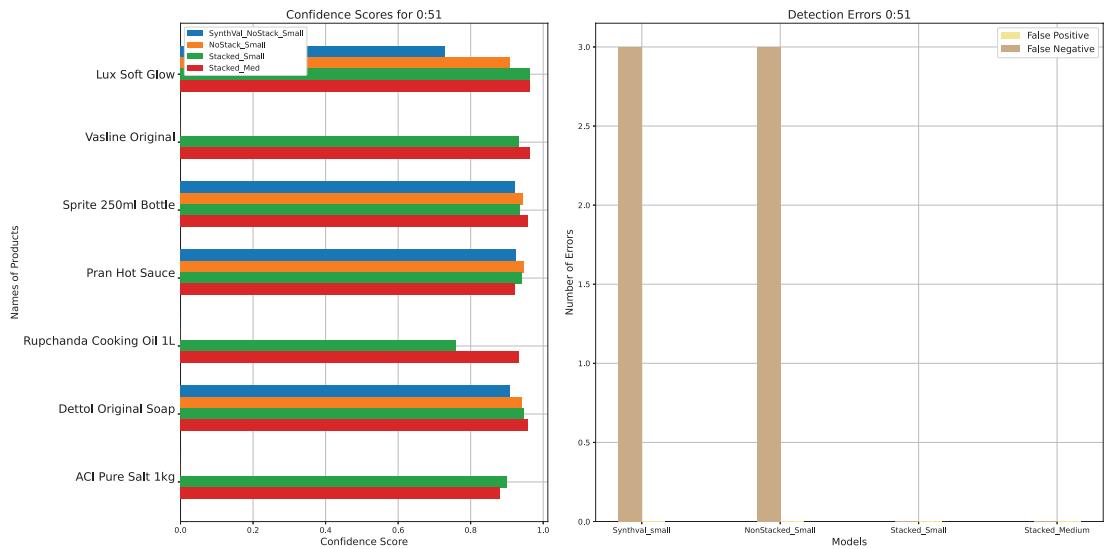


FIGURE 5.13: Analysis for time frame 0:51

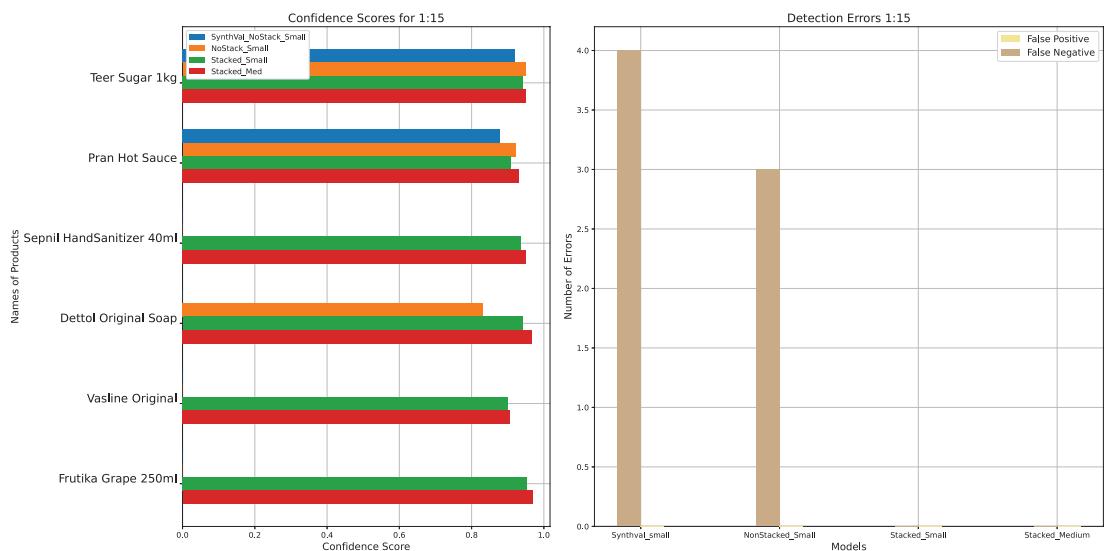


FIGURE 5.14: Analysis for time frame 1:15

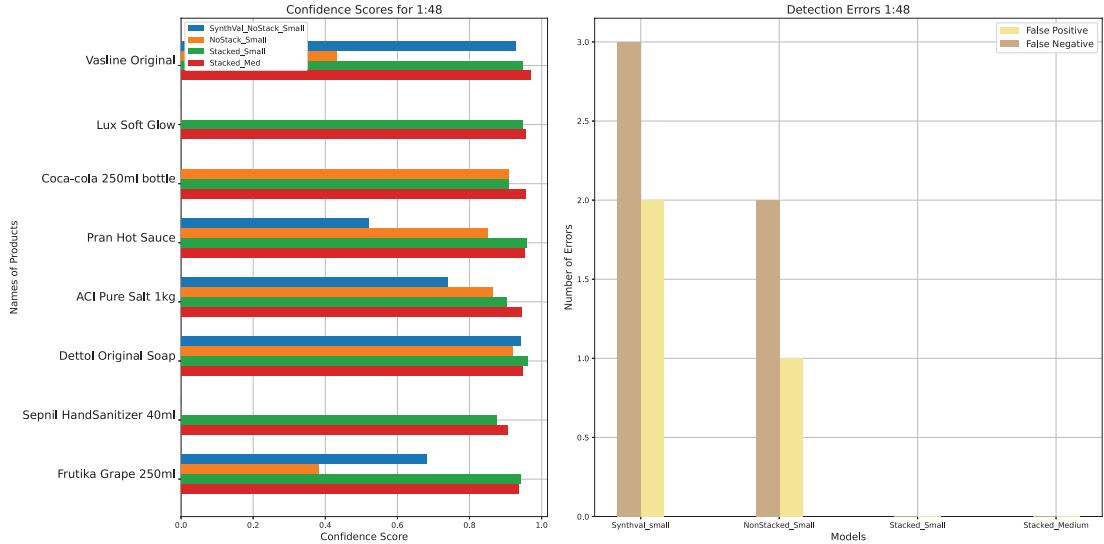


FIGURE 5.15: Analysis for time frame 1:48

From all the figures above, it is significantly noticeable that after adding the staked images to the dataset, real-time detection has improved. The confidence scores are consistent and there are zero detection errors in the selected time frames. The performance for both small and medium models are very similar. In a later section, real-time performance change due to hyperparameter evolution has been discussed.

5.9 Analysis on Different Backgrounds

The confusion matrices for the model's performances on Test Set - 3 are shown in this section.

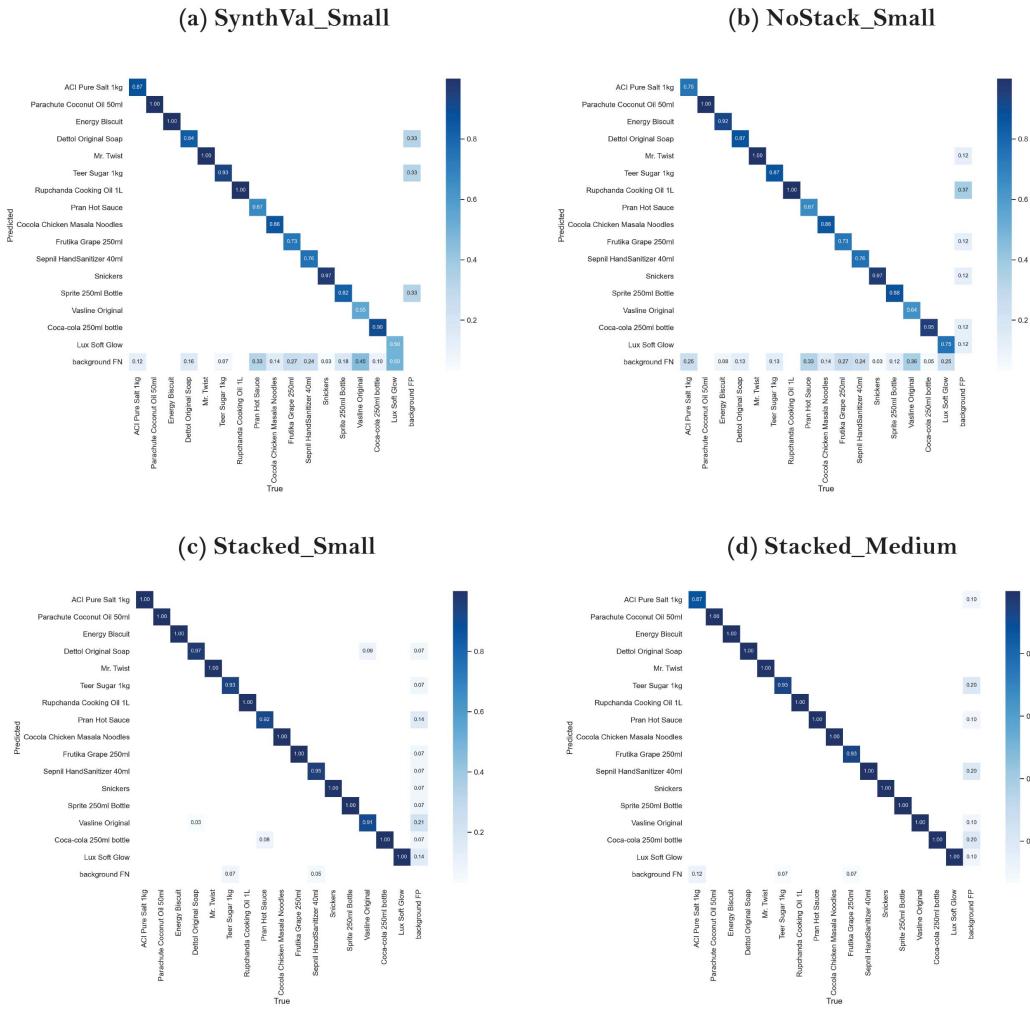


FIGURE 5.16: Confusion Matrices on Test Set - 3

From the matrices, we can notice significant improvement after adding densely arranged images for the model in (c) . Also, a small but noticeable improvement is also observed between (c) and (d) after shifting to the medium model which is consistent with the previous analyses.

5.10 Results after Data Augmentation

In Chapter 4, we have discussed the inclusion of synthetic data in our dataset. Here we compare the performance changes of the models based on mAP scores due to the addition of synthetic data. For simplicity, we have only compared the results of our final model (Stacked_Medium).

In the Figure 5.17, each group of 4 bars represent a test set. The first subgroup of 2 bars represent mAP@0.5 and the other subgroup represents mAP@0.5:0.95. The first bar and the second bar in each subgroup show performance without and with data augmentation respectively. It can be noticed that mAP scores have improved almost 10% in all cases due to augmentation via synthetic data.

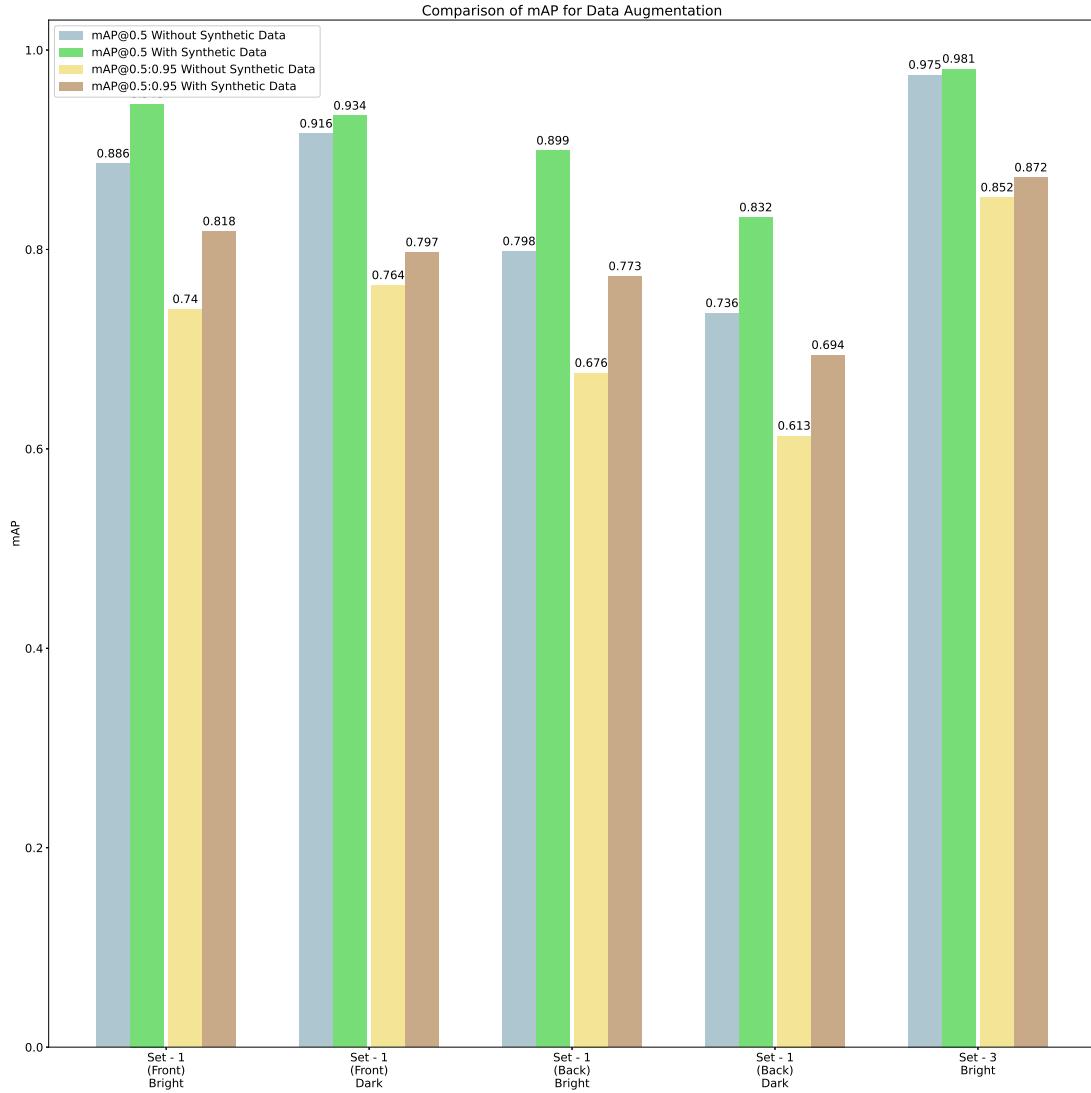


FIGURE 5.17: Comparison of mAP after Data Augmentation

5.11 Effects of Hyperparameter Evolution

In Chapter 4, we have also discussed hyperparameter evolution, here we discuss the results of the evolution and how it fairs against the default values we used before.

5.11.1 Values Obtained through evolution

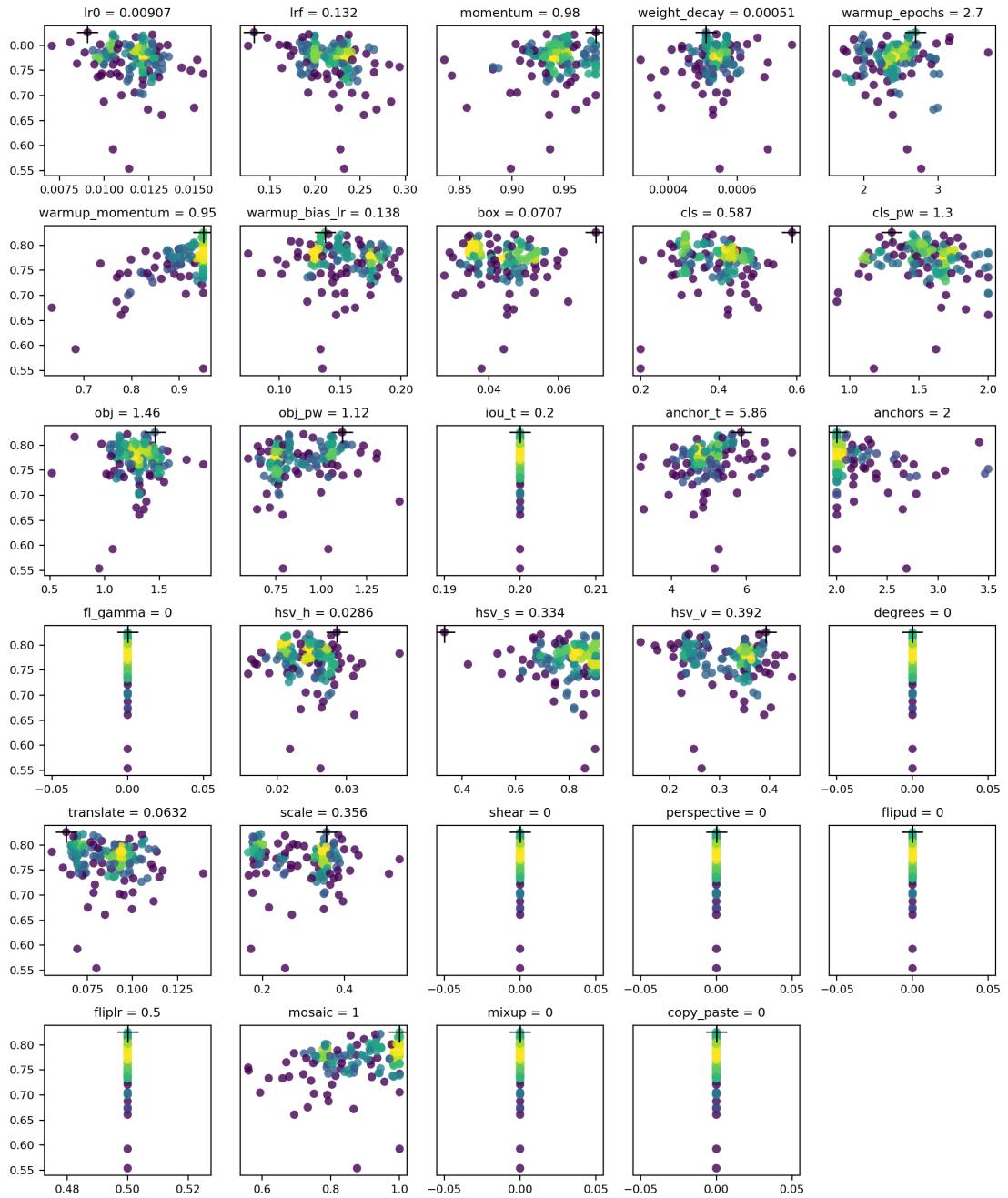


FIGURE 5.18: Distribution of Hyperparameter values for 150 generation of evolution

Based on the evolution, we obtained the best results on generation 11. Here are the hyperparameter values found in that generations:

lr0:	0.01269
lrf:	0.23471
momentum:	0.96013
weight_decay:	0.00035
warmup_epochs:	1.7472
warmup_momentum:	0.76925
warmup_bias_lr:	0.12307
box:	0.04871
cls:	0.38396
cls_pw:	2.0
obj:	1.17779
obj_pw:	1.03978
iou_t:	0.2
anchor_t:	6.14466
anchors:	2.72045
fl_gamma:	0.0
hsv_h:	0.02546
hsv_s:	0.83462
hsv_v:	0.44497
degrees:	0.0
translate:	0.11493
scale:	0.38548
shear:	0.0
perspective:	0.0
flipud:	0.0
fliplr:	0.5
mosaic:	0.91917
mixup:	0.0
copy_paste:	0.0

5.11.2 Performance comparison on test sets

We have trained the YOLOv5m model with both default and the generation 11 hyperparameters and tested them on our 3 test sets. The chart on Figure 5.19 is structured in the same ways as the previous one. We can observe that after hyperparameter evolution the performance has noticeably dropped on this test set. Each generation had been trained for 10 epochs, which may be the cause for these unsatisfactory results. 10 epochs of runtime cannot account for the scenario of 100 epochs, so the performance didn't scale proportionally. However, we do notice some improvements in real-time detection. On our Test Set -2, the model performs similarly well comparing to our model with default hyperparameters and improves on the consistency of detection. The detections for the evolved hyperparameters are much crisper.

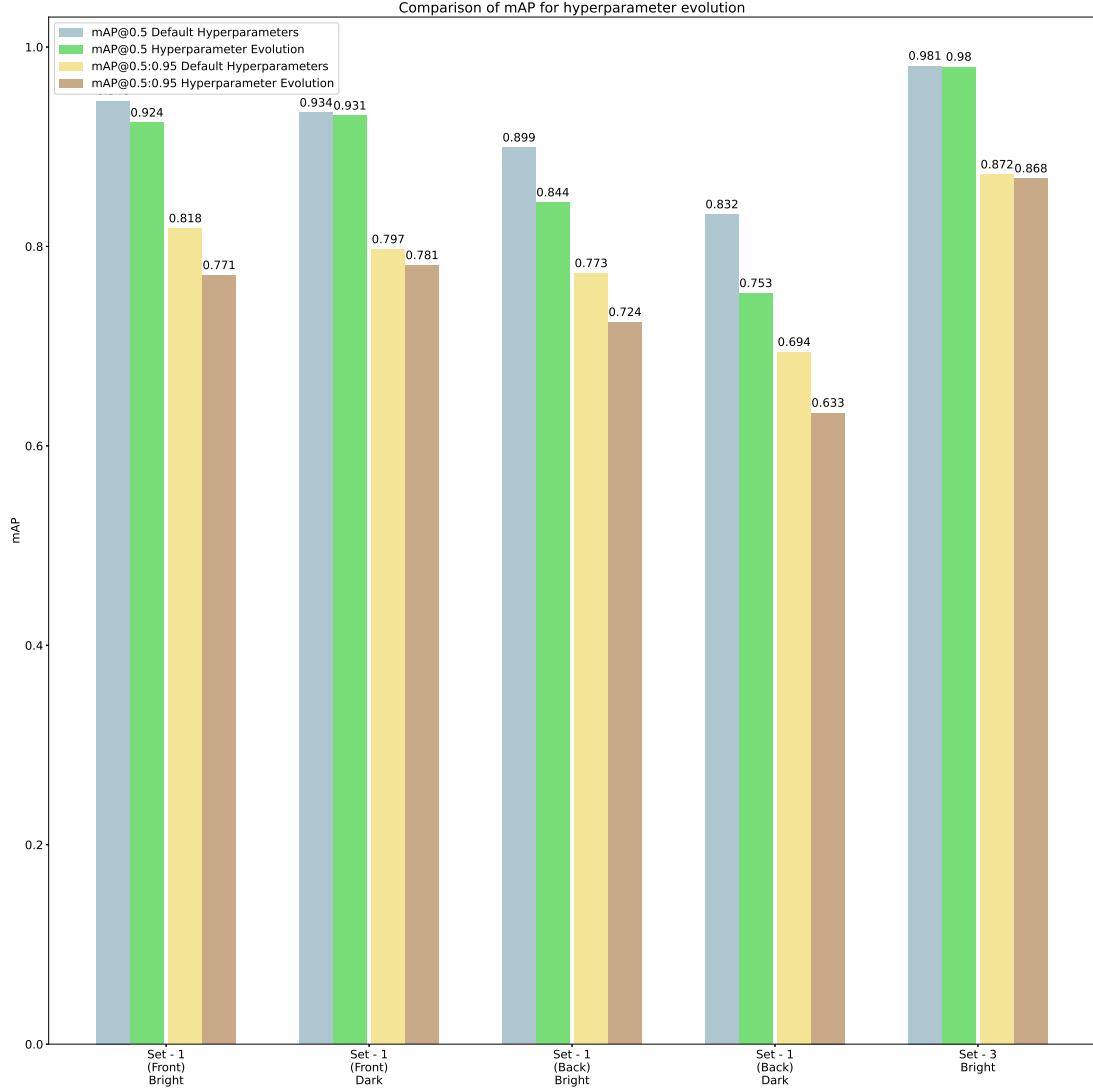


FIGURE 5.19: Comparison of mAP after hyperparameter evolution

In the Figure 5.20, we compare some frames of Test Set - 2. Frames 109-1019 are shown in one frame intervals. When the “Pran Hot Sauce” is placed into the scene the default hyperparameter model detects an additional “Rupchanda Cooking Oil” (lime green label) on frames 1013-1017 (marked with red circles). This is most likely due to the hand of the person and the moving bottle. But we notice no such glitches in the detection of the model after hyperparameter evolution.

Frame Anaylsis for Test Set 2 after Hyperparameter Evolution

Frame	1009	1011	1013
Default Hyperparameter			
After Hyperparameter Evolution			
Frame	1015	1017	1019
Default Hyperparameter			
After Hyperparameter Evolution			

FIGURE 5.20: Frame Analysis for Test Set 2 after Hyperparameter Evolution

CHAPTER 6

CONCLUSION AND FUTURE SCOPE

A CV-based billing system for retail products has been discussed in this project. The objective of the project was to replace traditional barcode scanners in supermarkets in order to speed up the checkout process. We have implemented a GUI-based checkout application that works on a YOLO-based object detection system to bill the products. We also analyze the efficiency of the YOLO-based models and how several factors like data augmentation and the increase of densely arranged images affect the performance.

6.1 Discussion

In Chapter 5, we analyzed our models on multiple test sets each focussing on different aspects of object detection. We have found that adding a few densely arranged product photos substantially increased our performance. It is apparent from the data augmentation also adds a significant improvement to the models. Furthermore, we found that using a larger model showed a very slight improvement over the small model. And in the case of hyperparameter tuning, we observed some improvement in the quality of real-time detection but showed a dip in performance on the other test sets. This is due to the fact that we only could train

150 generations and that up to 10 epochs which isn't sufficient to generate a hyperparameter-based performance improvement in the final model.

In the several tests we conducted, the final model (Stacked_Medium) performed reasonably well in both front and back-facing labels. The model also performs well independent of lighting conditions, only a small dip in performance is observed in darker environments. Back labels in the darker lighting conditions proved to be the most difficult for the models, but our final model performs reasonably well with mAP@0.5 and mAP@0.5:0.95 around 0.83 and 0.7.

6.2 Future Scopes

Apart from the results stated in our analysis, there are certain limitations to our system that can be improved in the future by applying different methods. The limitations and scopes for improvement are as follows:

1. As our approach towards detecting products is based on the shape and pattern of the labels, the model will fail to detect differences between products that have exactly similar-looking packaging but have different sizes. For example, Coca-cola bottles of 1L and 2L may not be detected as different products by our system. This may be solved by adding extra conditions for bounding box sizes and product orientation.
2. Single products cannot be removed from the list, this can be solved by improving the UI
3. Product labels that have small changes in them may not be efficiently detected as well. This can be circumvented by creating an ensemble of multiple models and use another model like one short learning method to determine the small differences [39].
4. Not all supermarket goods are packaged products, products like rice, meat which are often sold based on weight cannot be billed via this system. The implementation of a simultaneous QR Code based billing with image recognition [110] may be a good approach.

5. Even though we have created a CV-based solution, the system still requires a cashier. We can modify the system, incorporating a conveyor belt mechanism similar to the ARC Vision System [13] and develop an automated solution that will have all the advantages of our system.
6. Even though real-time detection performance may not be absolutely necessary for billing purposes as it is not essential to analyze every frame in a supermarket scenario. YOLO-based real-time detection works well on GPU-based systems which are expensive to deploy. More research is needed to optimize the model for CPU-based detection in order to make it cost-effective. One solution can be using ONNX [111] to convert to a more CPU-optimized framework like OpenVINO [112] will substantially improve performance on CPU-based systems. The inclusion of automated product type-based sorting and packaging can be a good future scope for the system, in order to move towards a fully automated supermarket experience.

BIBLIOGRAPHY

- [1] S. J. Pan and Q. Yang, “A survey on transfer learning,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, 2010.
- [2] Z. Zou, Z. Shi, Y. Guo, and J. Ye, “Object detection in 20 years: A survey,” 2019.
- [3] M. Rieder and R. Verbeet, “Robot-human-learning for robotic picking processes,” in *Hamburg International Conference of Logistics (HICL) 2019*. epubli GmbH, 2019, pp. 87–114.
- [4] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. [Online]. Available: https://www.cv-foundation.org/openaccess/content_cvpr_2016/html/Redmon_You_Only_Look_CVPR_2016_paper.html
- [5] M. Menegaz. Understanding yolo. [Online]. Available: <https://medium.com/hackernoon/understanding-yolo-f5a74bbc7967>
- [6] Cnn bounding box predictions. [Online]. Available: <http://datahacker.rs/deep-learning-bounding-boxes/>
- [7] L. Weng. Object detection part 4: Fast detection models. [Online]. Available: <https://lilianweng.github.io/lil-log/2018/12/27/object-detection-part-4.html>
- [8] R. Takahashi, T. Matsubara, and K. Uehara, “Data augmentation using random image cropping and patching for deep cnns,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 9, pp. 2917–2931, 2019.

- [9] S. I. Nikolenko *et al.*, “Synthetic data for deep learning,” *arXiv preprint arXiv:1909.11512*, vol. 3, p. 11, 2019.
- [10] Laughing-q. Yolov5 network structure. [Online]. Available: <https://blog.csdn.net/Q1u1NG/article/details/107511465>
- [11] R. Morabito and F. C. R. de Lima, “A markovian queueing model for the analysis of user waiting times in supermarket checkouts,” *International Journal of Operations and Quantitative Management*, vol. 10, no. 2, pp. 165–177, July 2004.
- [12] *Amazon Dash Cart*, Last accessed on October 2021. [Online]. Available: <https://www.amazon.com/b?ie=UTF8&node=21289116011>
- [13] S. T. Bukhari, A. W. Amin, M. A. Naveed, and M. R. Abbas, “Arc: A vision-based automatic retail checkout system,” 2021.
- [14] M. Zeiler, M. Ranzato, R. Monga, M. Mao, K. Yang, Q. Le, P. Nguyen, A. Senior, V. Vanhoucke, J. Dean, and G. Hinton, “On rectified linear units for speech processing,” in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2013, pp. 3517–3521.
- [15] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [16] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [17] B. Santra and D. P. Mukherjee, “A comprehensive survey on computer vision based approaches for automatic identification of products in retail store,” *Image and Vision Computing*, vol. 86, pp. 45–63, 2019. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0262885619300277>
- [18] A. Ray, N. Kumar, A. Shaw, and D. P. Mukherjee, “U-pc: Unsupervised planogram compliance,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.

- [19] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, “Speeded-up robust features (surf),” *Computer Vision and Image Understanding*, vol. 110, no. 3, pp. 346–359, 2008, similarity Matching in Computer Vision and Multimedia. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1077314207001555>
- [20] D. Lowe, “Object recognition from local scale-invariant features,” in *Proceedings of the Seventh IEEE International Conference on Computer Vision*, vol. 2, 1999, pp. 1150–1157 vol.2.
- [21] ——, “Distinctive image features from scale-invariant keypoints,” *International Journal of Computer Vision*, vol. 60, p. 91–110, November 2004.
- [22] K. Iwamoto, R. Mase, and T. Nomura, “Bright: A scalable and compact binary descriptor for low-latency and high accuracy object identification,” in *2013 IEEE International Conference on Image Processing*, 2013, pp. 2915–2919.
- [23] J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, and Y. Gong, “Locality-constrained linear coding for image classification,” in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2010, pp. 3360–3367.
- [24] R. Smith, “An overview of the tesseract ocr engine,” in *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, vol. 2, 2007, pp. 629–633.
- [25] D. J. C. MacKay, “Information-Based Objective Functions for Active Data Selection,” *Neural Computation*, vol. 4, no. 4, pp. 590–604, 07 1992. [Online]. Available: <https://doi.org/10.1162/neco.1992.4.4.590>
- [26] A. M. Treisman and G. Gelade, “A feature-integration theory of attention,” *Cognitive Psychology*, vol. 12, no. 1, pp. 97–136, 1980. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/0010028580900055>
- [27] E. Frontoni, M. Contigiani, and G. Ribighini, “A heuristic approach to evaluate occurrences of products for the planogram maintenance,” in *2014*

- IEEE/ASME 10th International Conference on Mechatronic and Embedded Systems and Applications (MESA)*, 2014, pp. 1–6.
- [28] A. Franco, D. Maltoni, and S. Papi, “Grocery product detection and recognition,” *Expert Systems with Applications*, vol. 81, pp. 163–176, 2017. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0957417417301227>
- [29] L. Karlinsky, J. Shtok, Y. Tzur, and A. Tzadok, “Fine-grained recognition of thousands of object categories with single-example training,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [30] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman, “Return of the devil in the details: Delving deep into convolutional nets,” in *Proceedings of the British Machine Vision Conference*. BMVA Press, 2014.
- [31] N. Zheng, G. Loizou, X. Jiang, X. Lan, and X. Li, “Computer vision and pattern recognition,” *International Journal of Computer Mathematics*, vol. 84, no. 9, pp. 1265–1266, 2007. [Online]. Available: <https://doi.org/10.1080/00207160701303912>
- [32] K. Georgiadis, F. Kalaganis, P. Migkotzidis, E. Chatzilari, S. Nikolopoulos, e. D. Kompatsiaris, Ioannis”, D. Giakoumis, M. Vincze, and A. Argyros, “A computer vision system supporting blind people - the supermarket case,” in *Computer Vision Systems*. Cham: Springer International Publishing, 2019, pp. 305–315.
- [33] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015. [Online]. Available: <http://arxiv.org/abs/1409.1556>
- [34] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.

- [35] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, “Inception-v4, inception-resnet and the impact of residual connections on learning,” ser. AAAI’17. AAAI Press, 2017, p. 4278–4284.
- [36] T. Winlock, E. Christiansen, and S. Belongie, “Toward real-time grocery detection for the visually impaired,” in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Workshops*, 2010, pp. 49–56.
- [37] B. D. Lucas and T. Kanade, “An iterative image registration technique with an application to stereo vision,” in *Proceedings of the 7th International Joint Conference on Artificial Intelligence - Volume 2*, ser. IJCAI’81. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1981, p. 674–679.
- [38] L. Barrington, T. K. Marks, J. Hui-wen Hsiao, and G. W. Cottrell, “Nimble: A kernel density model of saccade-based visual memory,” *Journal of Vision*, vol. 8, no. 14, pp. 17–17, 11 2008. [Online]. Available: <https://doi.org/10.1167/8.14.17>
- [39] W. Geng, F. Han, J. Lin, L. Zhu, J. Bai, S. Wang, L. He, Q. Xiao, and Z. Lai, “Fine-grained grocery product recognition by one-shot learning,” in *Proceedings of the 26th ACM International Conference on Multimedia*, ser. MM ’18. New York, NY, USA: Association for Computing Machinery, 2018, p. 1706–1714. [Online]. Available: <https://doi.org/10.1145/3240508.3240522>
- [40] J. Liu and Y. Liu, “Grasp recurring patterns from a single view,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2013.
- [41] M. A. Fischler and R. C. Bolles, “Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography,” *Commun. ACM*, vol. 24, no. 6, p. 381–395, Jun. 1981. [Online]. Available: <https://doi.org/10.1145/358669.358692>
- [42] M. Merler, C. Galleguillos, and S. Belongie, “Recognizing groceries in situ using in vitro training data,” in *2007 IEEE Conference on Computer Vision and Pattern Recognition*, 2007, pp. 1–8.

- [43] A. Tonioni and L. Di Stefano, “Product recognition in store shelves as a subgraph isomorphism problem,” in *Image Analysis and Processing - ICIAP 2017*, S. Battiato, G. Gallo, R. Schettini, and F. Stanco, Eds. Cham: Springer International Publishing, 2017, pp. 682–693.
- [44] A. Abdel-Hakim and A. Farag, “Csift: A sift descriptor with color invariant characteristics,” in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, vol. 2, 2006, pp. 1978–1983.
- [45] K. van de Sande, T. Gevers, and C. Snoek, “Evaluating color descriptors for object and scene recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 9, pp. 1582–1596, 2010.
- [46] S. Leutenegger, M. Chli, and R. Y. Siegwart, “Brisk: Binary robust invariant scalable keypoints,” in *2011 International Conference on Computer Vision*, 2011, pp. 2548–2555.
- [47] *C930E BUSINESS WEBCAM DATASHEET*, Last accessed on October 2021. [Online]. Available: https://www.logitech.com/content/dam/logitech/en_us/video-collaboration/pdf/c930e-datasheet.pdf
- [48] H.-I. Suk, “Chapter 1 - an introduction to neural networks and deep learning,” in *Deep Learning for Medical Image Analysis*, S. K. Zhou, H. Greenspan, and D. Shen, Eds. Academic Press, 2017, pp. 3–24. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/B978012810408800002X>
- [49] S. Agatonovic-Kustrin and R. Beresford, “Basic concepts of artificial neural network (ann) modeling and its application in pharmaceutical research,” *Journal of Pharmaceutical and Biomedical Analysis*, vol. 22, no. 5, pp. 717–727, 2000. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0731708599002721>
- [50] F. Rosenblatt, “The perceptron: A probabilistic model for information storage and organization in the brain,” *Psychological Review*, vol. 65, no. 6, p. 386–408, 1958. [Online]. Available: <https://psycnet.apa.org/doi/10.1037/h0042519>

- [51] P. Y. Simard, Y. A. LeCun, J. S. Denker, and B. Victorri, *Transformation Invariance in Pattern Recognition — Tangent Distance and Tangent Propagation*. Berlin, Heidelberg: Springer Berlin Heidelberg, 1998, pp. 239–274.
- [52] L. Torrey and J. Shavlik, *Transfer Learning*. IGI Global, 2010, pp. 242–264.
- [53] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- [54] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Commun. ACM*, vol. 60, no. 6, p. 84–90, May 2017. [Online]. Available: <https://doi.org/10.1145/3065386>
- [55] C. Szegedy, A. Toshev, and D. Erhan, “Deep neural networks for object detection,” in *Advances in Neural Information Processing Systems*, C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, Eds., vol. 26. Curran Associates, Inc., 2013. [Online]. Available: <https://proceedings.neurips.cc/paper/2013/file/f7cade80b7cc92b991cf4d2806d6bd78-Paper.pdf>
- [56] *A Gentle Introduction to Object Recognition With Deep Learning*, Last accessed on October 2021. [Online]. Available: <https://machinelearningmastery.com/object-recognition-with-deep-learning/>
- [57] R. Girshick, “Fast r-cnn,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.
- [58] K. He, G. Gkioxari, P. Dollar, and R. Girshick, “Mask r-cnn,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [59] *R-CNN – Neural Network for Object Detection and Semantic Segmentation*, Last accessed on October 2021. [Online]. Available: <https://neurohive.io/en/popular-networks/r-cnn/>
- [60] C. Chen, M.-Y. Liu, O. Tuzel, and J. Xiao, “R-cnn for small object detection,” in *Computer Vision – ACCV 2016*, S.-H. Lai, V. Lepetit, K. Nishino,

- and Y. Sato, Eds. Cham: Springer International Publishing, 2017, pp. 214–230.
- [61] S. Lazebnik, C. Schmid, and J. Ponce, “Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories,” in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, vol. 2, 2006, pp. 2169–2178.
- [62] K. He, X. Zhang, S. Ren, and J. Sun, “Spatial pyramid pooling in deep convolutional networks for visual recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [63] A. Bochkovskiy, C.-Y. Wang, and H.-y. Liao, “Yolov4: Optimal speed and accuracy of object detection,” 04 2020.
- [64] G. Maindola. A brief history of yolo object detection models from yolov1 to yolov5. [Online]. Available: https://machinelearningknowledge.ai/a-brief-history-of-yolo-object-detection-models/#YOLOv4_8211_Optimal_Speed_and_Accuracy_of_Object_Detection
- [65] Yolo: Real-time object detection. [Online]. Available: <https://pjreddie.com/darknet/yolo/>
- [66] R. Gandhi. R-cnn, fast r-cnn, faster r-cnn, yolo — object detection algorithms. [Online]. Available: <https://towardsdatascience.com/r-cnn-fast-r-cnn-faster-r-cnn-yolo-object-detection-algorithms-36d53571365e>
- [67] Overview of the yolo object detection algorithm. [Online]. Available: <https://medium.com/@ODSC/overview-of-the-yolo-object-detection-algorithm-7b52a745d3e0>
- [68] J. Deng, X. Xuan, W. Wang, Z. Li, H. Yao, and Z. Wang, “A review of research on object detection based on deep learning,” *Journal of Physics: Conference Series*, vol. 1684, p. 012028, nov 2020. [Online]. Available: <https://doi.org/10.1088/1742-6596/1684/1/012028>
- [69] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” 2016.

- [70] A. V. Thatte. Evolution of yolo — yolo version 1. [Online]. Available: <https://towardsdatascience.com/evolution-of-yolo-yolo-version-1-afb8af302bd2>
- [71] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *Proceedings of the 32nd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, F. Bach and D. Blei, Eds., vol. 37. Lille, France: PMLR, 07–09 Jul 2015, pp. 448–456. [Online]. Available: <https://proceedings.mlr.press/v37/ioffe15.html>
- [72] J. Redmon and A. Farhadi, “Yolo9000: Better, faster, stronger,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [73] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. Berg, and L. Fei-Fei, “Imagenet large scale visual recognition challenge,” *International Journal of Computer Vision*, vol. 115, 09 2014.
- [74] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [75] J. Redmon and A. Farhadi, “Yolov3: An incremental improvement,” 2018.
- [76] C.-Y. Wang, H.-Y. M. Liao, Y.-H. Wu, P.-Y. Chen, J.-W. Hsieh, and I.-H. Yeh, “CspNet: A new backbone that can enhance learning capability of cnn,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2020.
- [77] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, “Densely connected convolutional networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [78] J. Solawetz. Yolov5 new version - improvements and evaluation. [Online]. Available: <https://blog.roboflow.com/yolov5-improvements-and-evaluation/>

- [79] ——. How to train yolov5 on a custom dataset. [Online]. Available: <https://blog.roboflow.com/how-to-train-yolov5-on-a-custom-dataset/>
- [80] A. Mikołajczyk and M. Grochowski, “Data augmentation for improving deep learning in image classification problem,” in *2018 International Interdisciplinary PhD Workshop (IIPhDW)*, 2018, pp. 117–122.
- [81] C. Shorten and T. M. Khoshgoftaar, “A survey on image data augmentation for deep learning,” *Journal of Big Data*, vol. 6, no. 1, pp. 1–48, 2019.
- [82] I. Limit. Composing images with python for synthetic datasets. [Online]. Available: <https://www.immersivelimit.com/tutorials/composing-images-with-python-for-synthetic-datasets>
- [83] S. Pokhrel. Image data labelling and annotation-everything you need to know. [Online]. Available: <https://towardsdatascience.com/image-data-labelling-and-annotation-everything-you-need-to-know-86ede6c684b1>
- [84] D. Choi, C. J. Shallue, Z. Nado, J. Lee, C. J. Maddison, and G. E. Dahl, “On empirical comparisons of optimizers for deep learning,” *CoRR*, vol. abs/1910.05446, 2019. [Online]. Available: <http://arxiv.org/abs/1910.05446>
- [85] N. S. Keskar and R. Socher, “Improving generalization performance by switching from adam to SGD,” *CoRR*, vol. abs/1712.07628, 2017. [Online]. Available: <http://arxiv.org/abs/1712.07628>
- [86] N. Ketkar, “Stochastic gradient descent,” in *Deep learning with Python*. Springer, 2017, pp. 113–132.
- [87] P. Probst, M. N. Wright, and A.-L. Boulesteix, “Hyperparameters and tuning strategies for random forest,” *WIREs Data Mining and Knowledge Discovery*, vol. 9, no. 3, p. e1301, 2019. [Online]. Available: <https://wires.onlinelibrary.wiley.com/doi/abs/10.1002/widm.1301>
- [88] J. Brownlee. What is the difference between a parameter and a hyperparameter? [Online]. Available: <https://machinelearningmastery.com/difference-between-a-parameter-and-a-hyperparameter/>

- [89] G. Jocher. Hyperparameter evolution. [Online]. Available: <https://github.com/ultralytics/yolov5/issues/607>
- [90] J. H. Holland, “Genetic algorithms,” *Scientific American*, vol. 267, no. 1, pp. 66–73, 1992. [Online]. Available: <http://www.jstor.org/stable/24939139>
- [91] L. Tani, D. Rand, C. Veelken, and M. Kadastik, “Evolutionary algorithms for hyperparameter optimization in machine learning for application in high energy physics,” 11 2020.
- [92] *Daraz*, Last accessed on October 2021. [Online]. Available: <https://www.daraz.com.bd/>
- [93] *The Selenium Browser Automation Project*, Last accessed on October 2021. [Online]. Available: <https://www.selenium.dev/documentation/>
- [94] *Beautiful Soup Documentation*, Last accessed on October 2021. [Online]. Available: <https://www.crummy.com/software/BeautifulSoup/bs4/doc/>
- [95] *dupeGuru*, Last accessed on October 2021. [Online]. Available: <https://dupeguru.voltaicideas.net/>
- [96] *LabelImg*, Last accessed on October 2021. [Online]. Available: <https://github.com/tzutalin/labelImg>
- [97] R. Khandelwal. Coco and pascal voc data format for object detection. [Online]. Available: <https://towardsdatascience.com/coco-data-format-for-object-detection-a4c5eaf518c5>
- [98] A. Khan. Convert pascalvoc annotations to yolo. [Online]. Available: <https://gist.github.com/Amir22010/a99f18ca19112bc7db0872a36a03a1ec>
- [99] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran Associates, Inc.,

2019. [Online]. Available: <https://proceedings.neurips.cc/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf>
- [100] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Computer Vision – ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham: Springer International Publishing, 2014, pp. 740–755.
- [101] G. Jocher, A. Stoken, J. Borovec, NanoCode012, A. Chaurasia, TaoXie, L. Changyu, A. V, Laughing, tkianai, yxNONG, A. Hogan, lorenzomammana, AlexWang1900, J. Hajek, L. Diaconu, Marc, Y. Kwon, oleg, wanghaoyang0106, Y. Defretin, A. Lohia, ml5ah, B. Milanko, B. Fineran, D. Khromov, D. Yiwei, Doug, Durgesh, and F. Ingham, “ultralytics/yolov5: v5.0 - YOLOv5-P6 1280 models, AWS, Supervise.ly and YouTube integrations,” Apr. 2021. [Online]. Available: <https://doi.org/10.5281/zenodo.4679653>
- [102] W. Galitz, *The Essential Guide to User Interface Design: An Introduction to GUI Design Principles and Techniques*, ser. Wiley Desktop Editions. Wiley, 2007. [Online]. Available: https://books.google.com.bd/books?id=Q3Xp_Awu49sC
- [103] J. Willman, *Overview of PyQt5*. Berkeley, CA: Apress, 2021, pp. 1–42. [Online]. Available: https://doi.org/10.1007/978-1-4842-6603-8_1
- [104] *seaborn.pairplot*, Last accessed on October 2021. [Online]. Available: <https://seaborn.pydata.org/generated/seaborn.pairplot.html>
- [105] R. Padilla, S. L. Netto, and E. A. B. da Silva, “A survey on performance metrics for object-detection algorithms,” in *2020 International Conference on Systems, Signals and Image Processing (IWSSIP)*, 2020, pp. 237–242.
- [106] J. Brownlee. <https://machinelearningmastery.com/confusion-matrix-machine-learning/>. [Online]. Available: <https://machinelearningmastery.com/confusion-matrix-machine-learning/>

- [107] A. Rahmani, “Adapting google translate for english-persian cross-lingual information retrieval in medical domain,” in *2017 Artificial Intelligence and Signal Processing Conference (AISP)*, 2017, pp. 43–46.
- [108] *COCO Detection Evaluation*, Last accessed on October 2021. [Online]. Available: <https://cocodataset.org/#detection-eval>
- [109] *Open Images Challenge 2018 - object detection track - evaluation metric*, Last accessed on October 2021. [Online]. Available: https://storage.googleapis.com/openimages/web/object_detection_metric.html
- [110] Y. Gu and W. Zhang, “Qr code recognition based on image processing,” in *International Conference on Information Science and Technology*, 2011, pp. 733–736.
- [111] T. Le, G.-T. Bercea, T. Chen, A. Eichenberger, H. Imai, T. Jin, K. Kawachiya, Y. Negishi, and K. O’Brien, “Compiling onnx neural network models using mlir,” 08 2020.
- [112] Y. Gorbachev, M. Fedorov, I. Slavutin, A. Tugarev, M. Fatekhov, and Y. Tarkan, “Openvino deep learning workbench: Comprehensive analysis and tuning of neural networks inference,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, Oct 2019.

APPENDIX A: LIST OF ACRONYMS

COCO	Common Objects In Context
CNN	Convolutional Neural Network
CSPNet	Cross Stage Partial Network
CV	Computer Vision
DNN	Deep Neural Network
FPN	Feature Pyramid Network
GA	Genetic Algorithms
GAN	Generative Adversarial Networks
GUI	Graphical User Interface
FOV	Field of vision
OpenCV	Open Source Computer Vision Library
PANet	Path Aggregation Network
RANSAC	Random Sample Consensus
ReLU	Rectified Linear Unit
ROI	Region of Interest
SIFT	Scale Invariant Feature Transform
SGD	Stochastic Gradient Descent
SVM	Support Vector Machines
YOLO	You Only Look Once
IoU	Intersection over Union