



# **An Integrated System for Motorcycle Helmet and Overcrowding Violation Detection using Mask R-CNN with LLaVA-based Confirmation and License Plate Extraction**

A Report Submitted in

Partial Fulfillment of

requirement

for the Degree of

**Bachelor of Technology**

in

**Computer Science & Engineering**

By

**Harshika Singh (20214234)**

**Eleena Sarah Mathew (20214535)**

**Ayati Sonkar (20214300)**

**Bipul Singh Yadav (20214243)**

Under the guidance of

**Dr. Dushyant Kumar Singh**

**COMPUTER SCIENCE AND ENGINEERING DEPARTMENT**

**MOTILAL NEHRU NATIONAL INSTITUTE OF TECHNOLOGY ALLAHABAD  
PRAYAGRAJ**

# **UNDERTAKING**

We declare that the work presented in this report titled "**An Integrated System for Motorcycle Helmet and Overcrowding Violation Detection using Mask R-CNN with LLaVA-based Confirmation and License Plate Extraction**", submitted to the Computer Science and Engineering Department, Motilal Nehru National Institute of Technology Allahabad Prayagraj, for the award of the Bachelor of Technology degree in Computer Science & Engineering, is our original work. We have not plagiarized or submitted the same work for the award of any other degree. In case this undertaking is found incorrect, we accept that our degree may be unconditionally withdrawn.

Harshika Singh (20214234)

Eleena Mathew (20214535)

Ayati Sonkar (20214300)

Bipul Yadav (20214243)

# CERTIFICATE

I, Dr. Dushyant Kumar Singh, hereby certify that the work contained in the report titled "**An Integrated System for Motorcycle Helmet and Overcrowding Violation Detection using Mask R-CNN with LLaVA-based Confirmation and License Plate Extraction**" has been carried out under my supervision. I affirm that this work represents the original research efforts of the students and has not been submitted elsewhere for any degree or qualification.

Throughout the duration of this project, I have closely supervised the progress and methodology employed by the author. I can attest to the integrity and authenticity of the findings presented in this report.

Furthermore, I have ensured that all ethical considerations and academic standards have been diligently adhered to during the course of this research.

Should there be any inquiries or clarifications regarding the content or methodology of this report, I am available to provide further assistance and clarification.

Dr. Dushyant Kumar Singh

Associate Professor

Department of Computer Science & Engineering

Motilal Nehru National Institute of Technology Allahabad

6 May ,2025

## PREFACE

The purpose of this project is to explore the potential of using No Helmet and Overcrowding violation detection systems in enhancing road safety and traffic management. Specifically, we investigate the feasibility of employing advanced computer vision techniques and machine learning algorithms to detect and enforce traffic violations, such as helmet violations and overcrowding.

We undertook this project as a part of our Computer Science and Engineering program at Motilal Nehru National Institute of Technology, Allahabad, aiming to gain practical experience in the application of cutting-edge technologies in the transportation sector.

Our mentor, Dr. Dushyant Kumar Singh, played a pivotal role by providing invaluable guidance and support throughout the project's development. Their expertise and encouragement empowered us to navigate through the complexities of dataset preparation, model training, and system deployment with confidence and proficiency. We are extremely grateful for his support and motivation.

## **ACKNOWLEDGMENTS**

We would like to extend our heartfelt gratitude to Dr. Dushyant Kumar Singh for his unwavering commitment and exceptional support throughout every stage of this research project. Dr. Dushyant Kumar Singh's invaluable guidance, profound expertise, and steadfast mentorship have been instrumental in shaping the trajectory of our research endeavors. His insightful feedback, constructive criticism, and patient encouragement have continuously propelled us forward, transforming what initially seemed like a daunting task into a journey of growth and accomplishment.

Dr. Dushyant Kumar Singh's dedication to fostering academic excellence and his genuine investment in our academic and personal development have left an indelible mark on our scholarly pursuits. His mentorship transcends the boundaries of conventional academic supervision, serving as a beacon of inspiration and guidance that will continue to illuminate our paths long after the conclusion of this project.

In closing, we extend our profound gratitude to Dr. Dushyant Kumar Singh and all the esteemed panel members for their unwavering support, invaluable guidance, and profound contributions to our research journey. Their enduring impact will continue to resonate within our academic pursuits, inspiring us to strive for excellence and make meaningful contributions to our respective fields of study.

# **TABLE OF CONTENTS**

<b>Preface</b>	<b>4</b>
<b>Acknowledgement</b>	<b>5</b>
<b>Chapter1: Introduction</b>	<b>7</b>
<b>Chapter2: Proposed Work</b>	
<b>2.1 Scenario</b>	<b>9</b>
<b>2.2 Architecture</b>	<b>12</b>
<b>2.3 Implementation</b>	<b>17</b>
<b>Chapter3: Experimental Setup and Result Analysis</b>	
<b>3.1 Steps to setup Environment</b>	<b>39</b>
<b>3.2 Result Analysis</b>	<b>40</b>
<b>Chapter4: Conclusion and Future Work</b>	<b>41</b>
<b>4.1 Conclusion</b>	<b>41</b>
<b>4.2 Future Work</b>	<b>43</b>
<b>Bibliography</b>	<b>45</b>

# Chapter 1

## INTRODUCTION

Today, with the increasing number of vehicles on the road, transportation has changed significantly. While it has enabled greater mobility and convenience, it has also led to a surge in traffic violations like helmet non-compliance, overcrowding on two-wheelers etc. The importance of road transport, especially in densely populated regions, is overshadowed by the rising number of road injuries and fatalities due to such violations. With India's fast-growing vehicle population and urban expansion driven by economic growth, there is an urgent need to evaluate and upgrade the existing traffic surveillance and violation monitoring systems.

A promising frontier in this direction is the use of **multimodal AI agents** — intelligent systems capable of interpreting and analysing data from multiple sources such as video feeds, traffic sensor data and textual reports. By combining computer vision, natural language processing, and sensor analytics, these AI agents can provide a holistic understanding of traffic dynamics. For instance, while visual feeds help in detecting helmet and overcrowding violations, text-based reports from field officers and real-time traffic data can be analysed together to identify accident-prone zones, peak violation hours, and resource deployment strategies. Integrating multimodal AI systems in traffic analysis opens new avenues for proactive traffic management, better prediction of risk patterns, and dynamic law enforcement decisions.

In general, the rise in road accidents and traffic violations indicates a pressing requirement for fundamental changes in surveillance and enforcement methods. Automated systems leveraging cutting-edge technologies and data analytics offer a hopeful solution to improve road safety and safeguard human lives. It becomes evident, as we attempt to understand the complexities of contemporary road transportation, that the time has come for substantial advancements in traffic surveillance.

Limited enforcement resources combined with an ever-increasing number of vehicles make monitoring traffic violations a major challenge, particularly in developing nations. Current systems are labour-intensive, prone to human error, and often fail to capture critical details. Through advanced computer vision algorithms and AI-powered automation, developers can now build scalable, real-time systems for monitoring, flagging violations, and assisting authorities in enforcing road safety regulations effectively.

Two specific scenarios are addressed in this project:

1. Helmet Detection: Determine if two wheelers, and the passenger on pillion are wearing helmets required by traffic laws.
2. Overcrowding Detection: Detecting instability in case of more than the permissible number of passengers in two wheelers.

Additionally, these violations jeopardize safety and result in high number of injury and fatalities in road accidents. We automatically detect such violations, reducing human

dependency for both errors and violations, making it more accurate, and thus improving road safety.

## **Role of Deep Learning in Traffic Violation Detection**

Computer vision tasks have long been revolutionized by deep learning techniques, which enables the machines to learn complex patterns in images. The most popular methods are YOLO (You Only Look Once), SSD (Single Shot Multibox Detector), and Mask R-CNN, as each provides different methods in detection and segmentation. However, for tasks involving instance-level segmentation, like distinguishing between humans and helmets or identifying overcrowding, Mask R-CNN excels due to its:

- Pixel-Level Accuracy: By providing precise segmentation masks, the system is then able to reject further segmentation of closely packed objects.
- Multi-Class Detection: It can simultaneously detect multiple objects captured in a single frame (for example humans, helmets, motorcycles).
- Integration with Video Feeds: As it can process frames from video streams at its run-time, Mask R-CNN can readily detect the violations in real time.

## **Dataset and Tools**

To facilitate instance segmentation and bounding box generation of classes like person\_bike, human, motorcycle and helmet, the project makes use of datasets annotated using tools like Roboflow. By using labeled data, we train the model appropriately to spot violations in different environment settings such as different lighting, occlusions, and camera angles.

## **Objectives**

The primary objectives of this project are:

1. Designing a deep learning-based framework for helmet violation detection.
2. To ensure that the framework can achieve high accuracy on violation detection in real world scenarios.
3. That ensure traffic authorities to receive actionable insights which would then allow enforcement of traffic regulations efficiently.

# **Chapter 2**

## **PROPOSED WORK**

### **2.1 SCENARIO**

Due to the unique character of the two-wheeler infringement challenge, there is a critical need for such a system. Efforts to monitor and enforce these violations using conventional methods often fail because they are dynamic in nature. Utilizing technologies such as algorithms including Mask-RCNN. The project is an attempt to improve the precision and efficiency of traffic rule enforcement.

Unlike traditional approaches, this system not only points out violations but shows complete evidence of the violations. It records visual evidence of violations not only by flagging instances of non-compliance, but also images of the motorcycles that were violated and the faces of the riders on the motorcycles. This all-encompassing approach not only facilitates immediate enforcement actions, but it also yields a wealth of evidence for future retrospective reference, analysis, and in the unlikely event it ever becomes necessary, legal proceedings.

It utilizes instance segmentation to accurately determine violators and provide riders with their vehicles to enable accurate overcrowding detection. Visual evidence is captured by the system, including images of annotated violators, segmented helmets (or lack thereof) and overcrowded vehicles. This evidence can be kept and used upon immediate enforcement or for a future analysis and legal process store.

### **Why Mask R-CNN?**

When we dig into the specifics of this project, we realize that it is not about technological innovation just, but a safer, more accountable, and responsive road ecosystem. This system focuses not only on detection but also on robust evidence capture and strives to augment the overall effort to provide roadways that will prioritize the safety and well-being of all commuters.

Mask R-CNN is an advanced deep learning model which extends Faster R-CNN to include instance segmentation and is a highly appropriate choice when performing object detection and segmentation. As per this project where we are trying to detect helmets, motorcycles, people and overcrowding violations Mask R-CNN is best suited to it because of its robustness and versatility.

---

#### **1. Instance Segmentation Capability**

- Mask R-CNN gives us both bounding box detection and pixel level segmentation which is critical for detecting helmets as well as distinguishing objects that are close together (e.g bike, multiple riders).

- This guarantees precise object identification (e.g., a helmet on someone's head, someone riding on a bike), and as a result accurate violation detection.

## **2. Multi-Class Detection**

- The project requires detecting four distinct classes: It is person\_on\_bike, motorcycle, human and helmet.
- Unlike other architectures, architecture of Mask R-CNN is designed for multi-class object detection and can detect and segmenting multiple objects within a single image even in crowded traffic scenes.

## **3. Feature Pyramid Network (FPN) for Scale Invariance**

- Objects in traffic scenes appear of varying size (e.g. helmets are much smaller than motorcycles).
- Together with FPN, Mask R-CNN has high accuracy for detecting objects at any size or resolution.

## **4. Region Proposal Network (RPN)**

- To efficiently generate region proposals, Mask R-CNN uses an RPN to first discover objects and then refines those proposals to detect objects.
- The RPN speeds up detection and allows for attention to regions that are isotropically likely to contain objects of interest and reduces false positives and computational overhead.

## **5. Robust Against Occlusions**

- Traffic scenes usually have occlusions (riders partially occluding motorcycles or helmets).
- This enables Mask R-CNN to produce segmentation masks to separate out objects which overlap.

## **6. High Accuracy and Fine-Grained Detection**

- Fine grained analysis is possible by the model giving class probabilities, bounding box coordinates, and the segmentation masks.
- This accuracy is necessary to identify violations like overcrowding on two wheelers, where being close is not an inherent advantage, and ensure enforcement outcomes are based on that precision.

## **7. COCO Format Compatibility**

- Instance segmentation is supported by COCO format instance segmentation dataset, and Mask R-CNN supports this format.
- Exporting the annotated data into the COCO format and robustly integrating it with the model is made easy due to the use of the Roboflow.

## **Advantages Specific to This Project**

## **1. Helmet Detection**

- Because mask R-CNN can segment, it is able to distinguish helmets from other objects.
- It also aids in accurately flagging people without helmets, relative to models that use only bounding boxes.

## **2. Overcrowding Detection**

The model can find and count humans on a motorcycle so it can check for overcrowding violations.

- Precise identification of each individual is possible with segmentation masks for reliable overcrowding analysis.

## **3. Real-Time Detection in Complex Scenarios**

- The traffic environments are dynamic; there are diverse lighting conditions and object densities.
- Mask R-CNN has a very robust architecture coupled with the availability of robust pre trained weights and it also adapts well to these challenges to be deployed in real world.

## **4. Transfer Learning**

- Training is made faster and performance better using pre trained weights available on datasets such as COCO.
- This lowers the need of high dataset size; thereby leading the project to be cost and time efficient.

Feature	Mask R-CNN	YOLO	SSD
Instance Segmentation	Yes	No	No
Precision	High	Moderate	Moderate
Multi-Class Detection	Robust	Good	Good
Handling Occlusions	Excellent	Moderate	Moderate
Scalability	Supports scale-invariant features	Limited	Limited

**Table 1: Table comparing Mask R-CNN with YOLO and SSD**

YOLO and SSD do not perform as well as Mask R-CNN in terms of segmentation, precision and handling occlusions, so we choose Mask R-CNN for this project.

---

## Why Multimodal AI Agents and LLaVA Integration?

As part of advancing traffic surveillance beyond visual detection, this project also considers the integration of **multimodal AI agents**. These agents can simultaneously process information from different sources — such as video streams, sensor data, traffic reports, and textual records — to offer a more comprehensive understanding of traffic dynamics.

A promising multimodal AI system is **LLaVA (Large Language and Vision Assistant)**, which combines large language models (LLMs) with computer vision capabilities. In a traffic analysis context, LLaVA can:

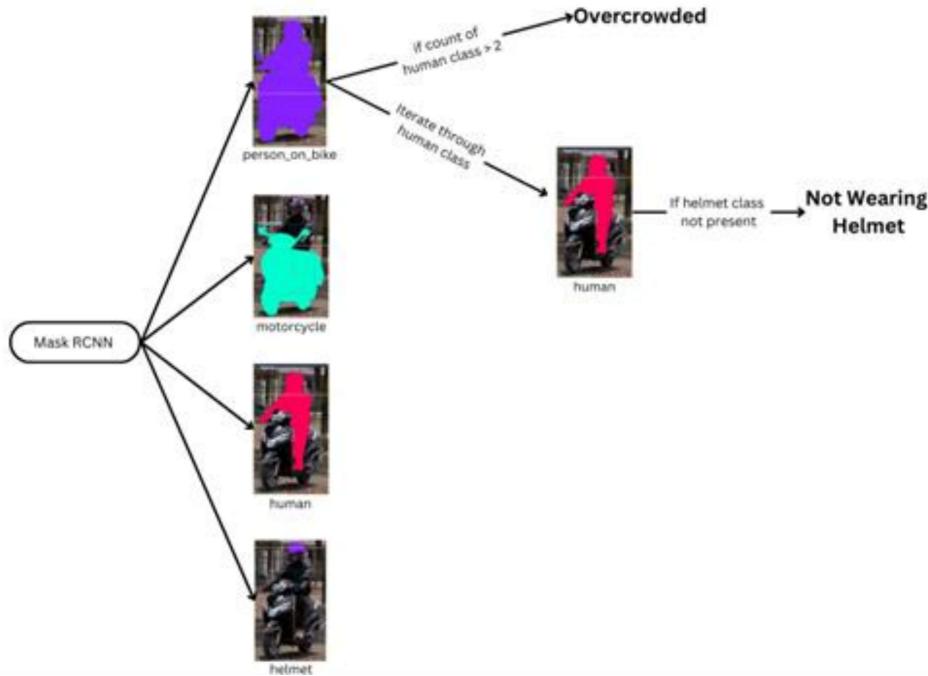
- Interpret textual traffic reports alongside visual evidence
- Generate automated violation summaries from detected images
- Cross-reference video feed detections with real-time road reports
- Provide contextual analysis about detected violations and suggest trends or enforcement strategies.

By incorporating LLaVA, the project could extend beyond raw detection to intelligent, interpretable, and context-aware traffic management, improving decision-making for authorities and reducing manual workloads.

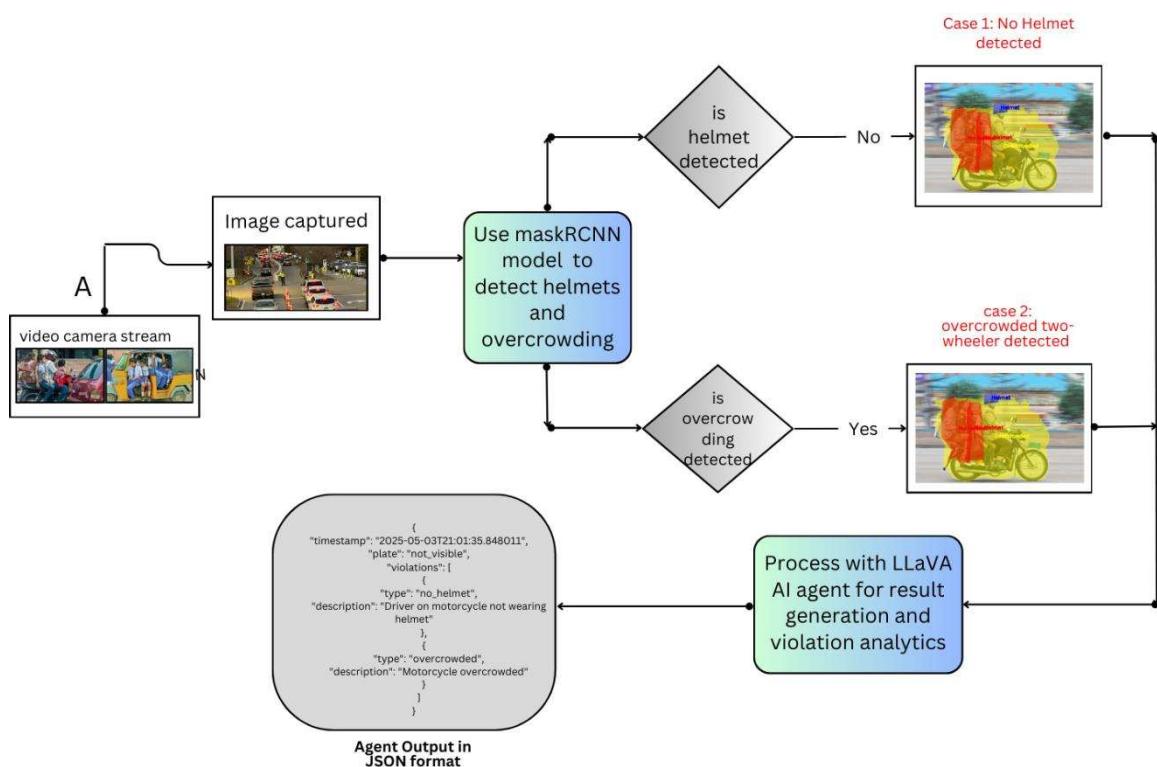
## Conclusion

Given the capabilities of **Mask R-CNN** in detecting, segmenting, and handling complex traffic scenes, it stands out as the best model for automated helmet violation and overcrowding detection. The proposed system, enhanced by potential multimodal AI agents like **LLaVA**, represents a significant advancement in traffic surveillance — creating a scalable, evidence-backed, real-time violation detection framework, designed to prioritize road safety and improve enforcement outcomes.

## 2.2 ARCHITECTURE



**Fig 1: Concise View of Feature Detection using Mask-RCNN**



**Fig 2: Block Diagram for Helmet and Overcrowding Detection using Mask-RCNN and LLaVA AI agent.**

The workflow for helmet violation and overcrowding detection on two wheelers using a Mask R-CNN model is represented in diagram (Fig 2).

It continuously captures live traffic footage on a video camera stream. It is monitoring two wheelers for helmet compliance and overcrowding violations. Video feed is extracted into frames. In fact, each frame serves as an input image for subsequent analysis. Such guarantees real-time detection and processing. Mask R-CNN model class then processes the captured image. It conducts instance segmentation and object detection and identifies core objects such as helmets, motorcycles, human and person\_on\_bike.

**In case 1:** If a motorcycle rider is wearing helmet or not is the analysis done by the model. If helmets are not detected, the system marks the detection as "Case 1: No Helmet".

**In case 2:** The model counts the number of Humans on a single motorcycle and checks if the count is more than the legal limit (raising the count to be greater than 2 humans for a two-wheeler) and raises a flag for 'Case 2: Overcrowded Two-Wheeler Detected' if the result is positive.

## Decision Nodes

- **Is Helmet Detected?**
  - **Yes** → No issue flagged.
  - **No** → Mark as **Case 1: No Helmet Detected**
    - Highlight the riders without helmets using a color-coded mask on the image.
- **Is Overcrowding Detected?**
  - **Yes** → Mark as **Case 2: Overcrowded Two-Wheeler Detected**
    - Highlight the overcrowded two-wheeler with segmented masks.

## Process with LLaVA AI Agent

- Both helmet and overcrowding violation images are processed by a **LLaVA (Large Language and Vision Assistant)** AI agent.
- This AI model interprets multimodal data (image annotations + possible metadata).
- Generates natural language descriptions of the violations, like:
  - "Driver on motorcycle not wearing helmet."
  - "Motorcycle overcrowded — three riders detected."

## Why LLaVA?

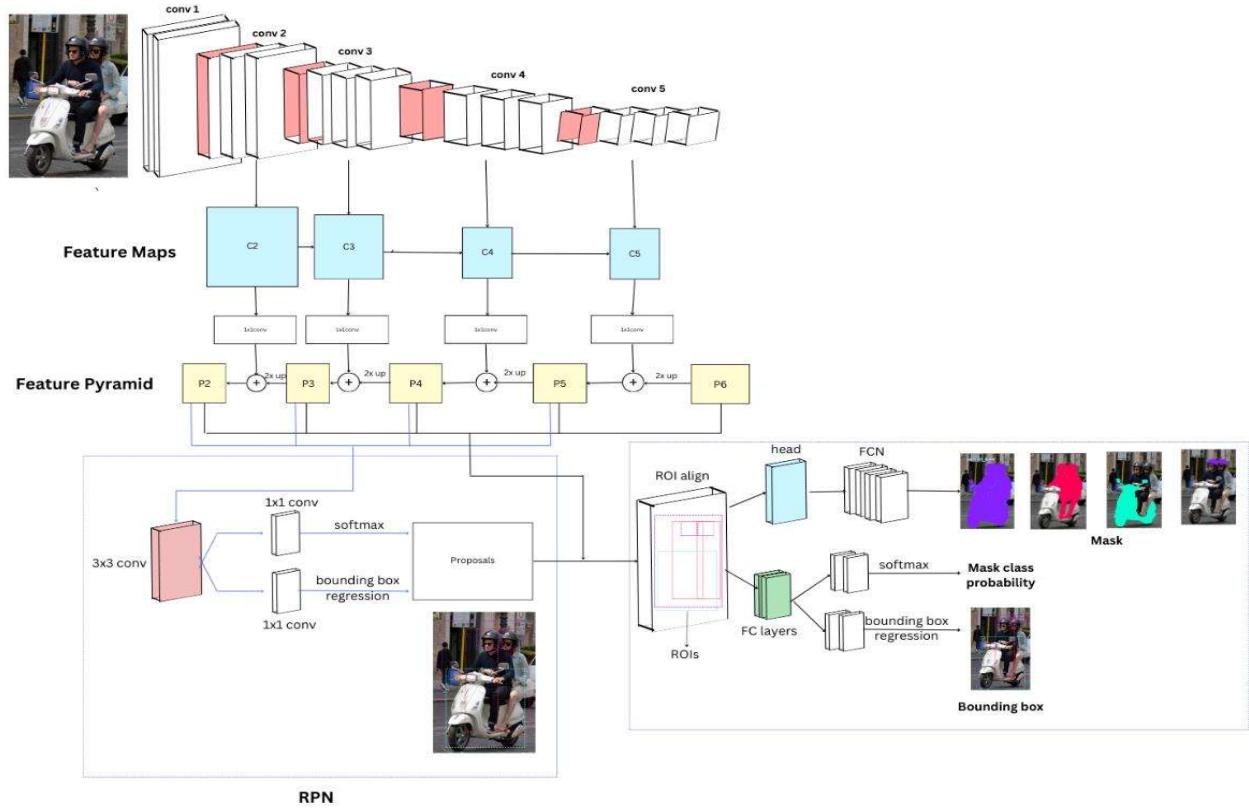
- It enhances Mask R-CNN's raw detection by:
  - Contextualizing violations.
  - Generating readable, structured outputs.
  - Possibly integrating location, license plate (if readable), and timestamp data.

## Agent Output in JSON Format

- LLaVA AI agent structures the output into a **JSON object**.
- JSON includes:
  - **timestamp**
  - **plate** (if available/visible)
  - **violations array**:
    - Type of violation (like "no\_helmet", "overcrowded").
    - Natural language description.

### Example Output:

```
{  
  "timestamp": "2025-05-03T12:03:15.840011",  
  "plate": "not_visible",  
  "violations": [  
    {  
      "type": "no_helmet",  
      "description": "Driver on motorcycle not wearing helmet"  
    },  
    {  
      "type": "overcrowded",  
      "description": "Motorcycle overcrowded"  
    }  
  ]  
}
```



**Fig 3: Mask R-CNN-Based Framework for Helmet Detection and Overcrowding Assessment on Two-Wheelers**

Video streams captured from surveillance cameras installed at critical traffic points such as the intersections and highways are processed with the help of the system (Fig3). In the real time the video streams are processed, individual frames are extracted and are analyzed further. A convolutional neural network (ResNet) is made to take in the input image.

Initially, layers conv1 to conv5 progressively extract the feature maps, which represent different spatial hierarchies and levels of abstraction up to the feature map. The outputs of the convolutional layers (C2, C3, C4, C5) are up sampled and combined together using lateral connection to form a feature pyramid (P2, P3, P4, P5, P6). This allows the system to work with objects of different scales through features analyzed out in different scales.

The **RPN** generates candidate object proposals:

1. **3x3 Convolution:** It scans the feature map to predict regions of interest (ROI).
2. **Bounding Box Regression and Object Classification:** It determines possible bounding boxes for objects (i.e., helmets, motorcycles, humans) and guesses the level of their background or foregroundness.
3. The final proposals are passed to the next stage where top proposals are selected and retained.

ROI Align is used to map proposals to the corresponding regions in the feature pyramid for spatial accuracy. It improves over ROI Pooling by addressing the quantizations errors.

Each ROI is processed by specialized sub-networks for:

- **Classification:** It classifies the object (rider, helmet or background etc).
- **Bounding Box Refinement:** Refines the bounding box coordinate precision.
- **Mask Prediction:** Fully Convolutional Network (FCN) based method for producing pixel wise masks for instance segmentation.

The segmentation masks that the mask layer produces indicate members of the group to the region — for example, riders or helmets. Objects are located and then highlighted with accurate bounding boxes. Outputs of the system are class labels and probabilities for each object, bounding boxes around detected objects, as well as instance segmentation masks for riders, helmets, and vehicles.

## **2.3 IMPLEMENTATION**

### **2.3.0 DATASET PREPARATION**

#### **1. Dataset Collection:**

- Images of motorcycle traffic scenes were collected, in which traffic also included motorcycles, riders, helmets, and people.
- From diverse scenarios, lighting conditions and viewpoints, the dataset is designed to be robust trained.

#### **2. Roboflow for Image Annotation:**

- Class Labels: \$our classes: person\_on\_bike, human, motorcycle, and helmet.
- Each object instance in the images was annotated with both:
  - Bounding Boxes: Regions labeled rectangular, and indicating the object's location.
  - Instance Segmentation Masks: Object boundaries marked precisely in pixel wise annotations.

#### **3. Augmentation in Roboflow:**

- Applied various data augmentation techniques to increase the dataset size and variability, such as:
  - Random cropping.
  - Horizontal (or) vertical flipping.
  - Brightness, contrast, and hue adjustments.
  - Scaling and rotation.

#### **4. Dataset Export:**

- The annotated dataset was exported from Roboflow in COCO format, which is compatible with Mask R-CNN.
- The COCO format includes:
  - Image metadata (size, file path).
  - Class labels plus bounding box coordinates.
  - Instance segmentation masks represented as polygon data.

## 5. Splitting the Dataset:

- For instance, we split the dataset into three sets; training, validation, test (e.g, 80% for training, 10% for validation, 10% for test).
- With balanced representation of all classes in each subset, biased results from one of the datasets were avoided.

## 6. Preprocessing for Mask R-CNN:

- The aspect ratio of all images was maintained (e.g., 640x640) by simply put zero padding wherever necessary.
- Then images and annotations were converted into a format that the Mask R-CNN model would be able to handle.

## 7. Dataset Validation:

- Consistency and correctness verified annotations (bounding boxes and masks).
- Found missing or overlapping annotations and cleaned the data for accuracy.

## 8. Final Output:

- The trained Mask R-CNN model for detecting and segmenting the objects against the four defined classes was trained using the processed dataset.

### 2.3.1 CONFIGURATIONS USED TO TRAIN THE MODEL

#### Model Architecture Parameters

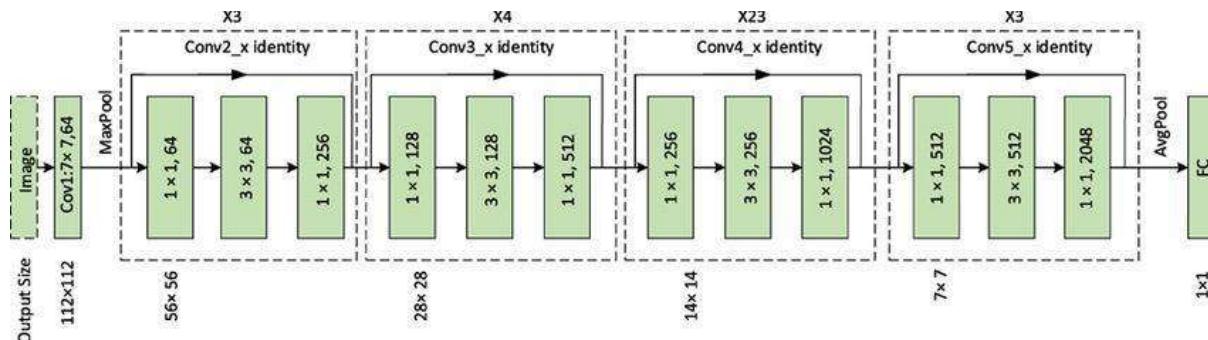
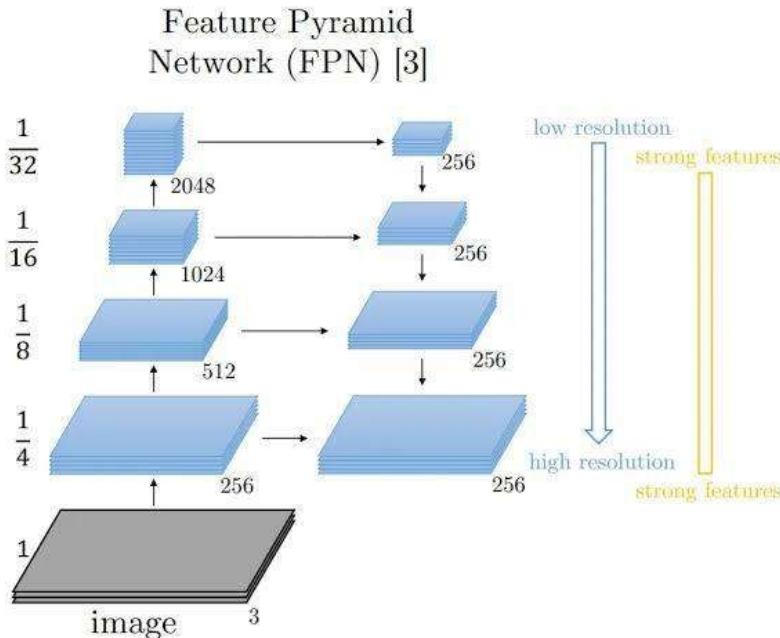


Fig 4: Resnet-101 Architecture

- **BACKBONE (resnet101)**: The feature extraction backbone used by the model, here ResNet-101, is a deep residual network for feature extraction.
- **BACKBONE\_STRIDES ([4, 8, 16, 32, 64])**: Strides of the backbone layers at different levels in the feature pyramid network (FPN).
- **FPN\_CLASSIF\_FC\_LAYERS\_SIZE (1024)**: Fully connected layer size for the classification head in FPN.



**Fig 5: Feature Pyramid Network**

## Image Properties

- **IMAGE\_MIN\_DIM (1024) / IMAGE\_MAX\_DIM (1024)**: The minimum and maximum dimensions of the input images after resizing.
- **IMAGE\_RESIZE\_MODE (square)**: Resizes images to a square while maintaining aspect ratio and adding padding as needed.
- **IMAGE\_SHAPE ([1024, 1024, 3])**: The final input shape of the images, including height, width, and number of channels (3 for RGB).
- **MEAN\_PIXEL ([123.7, 116.8, 103.9])**: Pixel mean values for normalization (used to center the pixel values of images).
- **IMAGE\_CHANNEL\_COUNT (3)**: Number of channels in the input image (3 for RGB).

---

## Training and Inference Parameters

- **BATCH\_SIZE (2)**: Number of images processed together per batch.
  - **IMAGES\_PER\_GPU (2)**: Number of images processed on a single GPU (should match BATCH\_SIZE for single-GPU setups).
  - **GPU\_COUNT (1)**: Number of GPUs being used for training.
  - **LEARNING\_RATE (0.001)**: Learning rate for the optimizer.
  - **LEARNING\_MOMENTUM (0.9)**: Momentum value for the optimizer.
  - **WEIGHT\_DECAY (0.0001)**: Regularization parameter to prevent overfitting.
  - **GRADIENT\_CLIP\_NORM (5.0)**: Clips gradients to prevent exploding gradients.
  - **STEPS\_PER\_EPOCH (100)**: Number of iterations per epoch.
  - **VALIDATION\_STEPS (10)**: Steps run during validation.
- 

## Object Detection Parameters

- **NUM\_CLASSES (3)**: Number of classes, including the background class. For example, if detecting "marble" and "granite", this would include 1 background class.
  - **DETECTION\_MIN\_CONFIDENCE (0.7)**: Minimum confidence threshold for a detection to be considered valid.
  - **DETECTION\_MAX\_INSTANCES (35)**: Maximum number of instances detected per image.
  - **DETECTION\_NMS\_THRESHOLD (0.3)**: Non-maximum suppression (NMS) threshold to remove duplicate detections.
- 

## Region Proposal Network (RPN)

- **RPN\_ANCHOR RATIOS ([0.5, 1, 2])**: Ratios of width to height for RPN anchors.
  - **RPN\_ANCHOR\_SCALES ((32, 64, 128, 256, 512))**: Anchor scales used by the RPN to generate anchors of varying sizes.
  - **RPN\_ANCHOR\_STRIDE (1)**: Stride of the anchors in the feature map.
  - **RPN\_TRAIN\_ANCHORS\_PER\_IMAGE (256)**: Number of anchors sampled for training.
  - **RPN\_NMS\_THRESHOLD (0.7)**: NMS threshold for filtering RPN proposals.
  - **PRE\_NMS\_LIMIT (6000)**: Maximum number of RPN proposals pre-NMS.
  - **POST\_NMS\_ROIS\_TRAINING (2000) / POST\_NMS\_ROIS\_INFERENCE (1000)**: Number of proposals kept post-NMS for training and inference.
- 

## RoI Pooling and Mask Generation

- **POOL\_SIZE (7)**: Size of the pooling layer for RoI Align for bounding box classification and regression.

- **MASK\_POOL\_SIZE (14)**: Pool size for mask segmentation.
  - **MASK\_SHAPE ([28, 28])**: Final shape of the masks.
  - **ROI\_POSITIVE\_RATIO (0.33)**: Fraction of positive RoIs used during training.
  - **TRAIN\_ROIS\_PER\_IMAGE (200)**: Number of RoIs per image used during training.
- 

### Mini Masks

- **USE\_MINI\_MASK(False)**: Disables using smaller masks for training. Useful when working with high-resolution images.
  - **MINI\_MASK\_SHAPE ((56, 56))**: Size of mini masks (not used since USE\_MINI\_MASK is False).
- 

### Loss Weights

- **LOSS\_WEIGHTS**: Defines the relative importance of different loss components:
    - rpn\_class\_loss: Loss for RPN classification.
    - rpn\_bbox\_loss: Loss for RPN bounding box refinement.
    - mrcnn\_class\_loss: Loss for Mask R-CNN classifier.
    - mrcnn\_bbox\_loss: Loss for bounding box refinement in Mask R-CNN.
    - mrcnn\_mask\_loss: Loss for mask prediction.
- 

These configurations have been fine-tuned for a specific dataset ("marble\_cfg\_coco") with 4 classes, aiming for accurate detection and segmentation of objects.

## 2.3.2 HELMET VIOLATION DETECTION

### 1. Image Upload:

- Traffic scene images are processed locally by users and used to provide traffic scene images to be analyzed.
- The uploaded image has to be exactly 640x640, and if not, then the image is resized by scaling and padding with zeros for the model to be compatible
- The system supports multiple image formats such as JPEG and PNG for seamless processing.

### 2. Image Preprocessing:

- Data preprocessing step is taken when the uploaded images are resized and normalized for model inference.

- It preprocesses the input so that the Mask R CNN model has uniform, high quality input.

### 3. Intersection over Union (IoU) Calculation

- To determine the spatial overlap between detected objects (e.g., between a human and a helmet), the system uses the Intersection over Union (IoU) metric. IoU calculates the ratio of the area of overlap to the area of union between two bounding boxes. This value helps decide if a detected human is wearing a helmet based on the overlap of their bounding boxes.

The formula used is:

$$\text{IoU} = \frac{\text{Area of Intersection}}{\text{Area of Union}}$$

```
def iou(box1, box2):
    y1 = max(box1[0], box2[0])
    x1 = max(box1[1], box2[1])
    y2 = min(box1[2], box2[2])
    x2 = min(box1[3], box2[3])

    inter_area = max(0, y2 - y1) * max(0, x2 - x1)
    box1_area = (box1[2] - box1[0]) * (box1[3] - box1[1])
    box2_area = (box2[2] - box2[0]) * (box2[3] - box2[1])

    union_area = min(box1_area, box2_area)
    return inter_area / union_area if union_area > 0 else 0
```

- A threshold of 0.5 was set to mark a person as wearing a helmet.

### 4. Helmet Violation Detection:

- The inputs for this go through the same preprocessing and the model performs instance segmentation as well as classifying objects like humans, helmets and motorcycles.
- The segmented objects using bounding boxes and masks are categorized into meaningful classes (motorcycle, person\_on\_bike, human, and helmet) and then detected.
- The system uses the Mask R-CNN model to detect helmets (class ID 1) in each image frame. Upon detection:
  - A colored mask is applied to the helmet region for visual marking.
  - The centroid of the detected helmet's bounding box is calculated.
  - A label “Helmet” is overlaid at the centroid location.

```

if class_id == 1:
    mask = masks[:, :, i]
    color = (255, 0, 0)
    image = apply_mask(image, mask, color)
    cx, cy = (x1 + x2) // 2, (y1 + y2) // 2
    cv2.putText(image, "Helmet", (cx, cy), cv2.FONT_HERSHEY_SIMPLEX,
                0.5, color, 2)

```

- This step ensures that helmets are properly highlighted in the processed images, aiding both system accuracy and result validation.
- Then again system employs the Mask R-CNN model to detect multiple object classes within a traffic scene, notably humans and helmets. After detection:
  - Each human detection (class ID 2) is checked for overlapping helmet detections using the Intersection over Union (IoU) metric.
  - If no helmet detection overlaps more than **50%** (IoU > 0.5) with a given human detection, the system considers that person as **not wearing a helmet**.

```

# Check for humans without helmets

if class_id == 2: # Human

    has_helmet = any(iou(boxes[i], helmet_box) > 0.5 for helmet_box in helmet_boxes)

    if not has_helmet:

        color = (0, 0, 255)

        mask = masks[:, :, i]

```

```

image = apply_mask(image, mask, color)

cx, cy = (x1 + x2) // 2, (y1 + y2) // 2 # Centroid of the bounding box

cv2.putText(image, "No Helmet", (cx, cy), cv2.FONT_HERSHEY_SIMPLEX, 0.5,
color, 2)

data['helmetViolation'] = False

```

- This visual marking and data flagging process ensures both visual and structured reporting of helmet violations in real-time.

## **5. Results and Visualization:**

- Visuals of the processed images with detected violations appear marked directly on top, with bounding boxes and segmentation masks imposed over those without helmets.
- Images are saved locally by the system for investigation of violations.

## **6. Scalable Deployment:**

This model and its processing pipeline are local and will perform reliably and quickly for helmet violation analysis.

## **7. Output:**

### **1. Detection Overview:**

The system has detected:

- **person\_bike** with a confidence score of **0.998**.

### **Visual Results:**

The following image illustrates the detection:



**Fig 6: Sample result for Image with helmet and No helmet and person\_bike**

#### **System Details:**

- **Input Image Shape:** (640, 640, 3)
- **Processed Image Shape:** (1024, 1024, 3)
- **Anchor Points:** 261,888
- **Confidence Scores:** person\_bike (0.998)

*Description:* With the detection system, the identified areas or the bike and the helmet, are effectively masked and the No helmet violation is detected.

#### **2. Detection Overview:**

The system has detected:

- **No Helmet for person\_bike.**
- Hence, No Helmet violation detected.

#### **Visual Results:**

The following image illustrates the detection:



**Fig 7: Sample result for Image with No helmet violation**

#### **System Details:**

- **Input Image Shape:** (640, 640, 3)
  - **Min Value:** 0.0
  - **Max Value:** 255.0
  - **Data Type:** uint8
- **Processed Image Shape:** (1, 1024, 1024, 3)
  - **Min Value:** -123.7
  - **Max Value:** 151.1
  - **Data Type:** float64
- **Image Metadata:** (1, 15)
  - **Min Value:** 0.0
  - **Max Value:** 1024.0
  - **Data Type:** float64
- **Anchor Points:** (1, 241, 778, 4)
  - **Min Value:** -0.3539
  - **Max Value:** 1.29134
  - **Data Type:** float32

#### **Description:**

The detection system takes the input image, processes it and normalize it into the specified range and use anchor points for object detection. The areas which are identified are masked.

### **2.3.3 Overcrowding Detection**

#### **1. Image Upload:**

- Traffic scene images are processed locally by users and used to provide traffic scene images to be analyzed.
- The uploaded image has to be exactly 640x640, and if not, then the image is resized by scaling and padding with zeros for the model to be compatible
- The system supports multiple image formats such as JPEG and PNG for seamless processing.

## **2. Image Preprocessing:**

- Data preprocessing step is taken when the uploaded images are resized and normalized for model inference.
- It preprocesses the input so that the Mask R CNN model has uniform, high quality input.

## **3. Box Inclusion Check (is\_inside Function)**

- To identify overcrowding, it is necessary to check whether multiple human detections are within the same motorcycle bounding box. For this, a utility function `is_inside` is used, which checks whether one bounding box is completely enclosed within another.

### **Algorithm Implementation:**

```
def is_inside(box1, box2):
    x1, y1, x2, y2 = box1
    x3, y3, x4, y4 = box2
    return (x1 >= x3 and y1 >= y3 and x2 <= x4 and y2 <= y4)
```

- This function helps count how many humans are present on a motorcycle and raises an overcrowding violation if more than two human detections are enclosed within a single motorcycle detection.

## **4. Overcrowding Violation Detection:**

In addition to helmet violations, the system identifies instances of motorcycle overcrowding. The Mask R-CNN model detects the `person_on_bike` class (class ID 4) in each image frame. The following approach is used:

- For every detected *person\_on\_bike* instance, the system counts the number of human detections overlapping with it.
- Intersection over Union (IoU) is calculated between each detected human and the *person\_on\_bike* bounding box. An IoU greater than **0.25**, or if a human box is entirely inside the *person\_on\_bike* box, is considered a valid overlap.
- If the number of overlapping humans exceeds **two**, the instance is marked as **overcrowded**.
- A mask is applied over the overcrowded instance, and a label “*Overcrowded*” is displayed at the centroid of the bounding box.
- **The implementation of this logic is shown below:**

```

if class_id == 4:

    # Count number of human boxes overlapping with this motorcycle box (IoU > 0.25) or
    # completely inside it

    human_count = sum( iou(boxes[i], human_box) > 0.25 or is_inside(human_box,
                           boxes[i]) for human_box in human_boxes

    )

    if human_count > 2: # Overcrowded if more than 2 humans detected
        color = (0, 255, 255)
        mask = masks[:, :, i]
        image = apply_mask(image, mask, color)

        # Calculate center of the motorcycle's bounding box
        cx, cy = (x1 + x2) // 2, (y1 + y2) // 2
        cv2.putText(image, "Overcrowded", (cx, cy),
                    cv2.FONT_HERSHEY_SIMPLEX, 0.5, color, 2)

        # Flag violation in the data dictionary
        data['overcrowdedViolation'] = True

```

- This method ensures accurate detection of motorcycle overcrowding, contributing to effective traffic rule enforcement.

## 5. Spatial Relationship Functions

To assess the positional relationships between detected objects within the image, two helper functions were implemented:

- **Intersection over Union (IoU):**  
This function calculates the ratio of the overlapping area between two bounding boxes to their union area. It is widely used to quantify the degree of overlap between two regions.

- **Is Inside:**

This function verifies whether one bounding box is completely contained within another, which is particularly useful in determining if a detected human lies entirely within a *person\_on\_bike* detection.

**The implementations of these functions are shown below:**

```
def is_inside(outer_box, inner_box): """ Check if one bounding box (inner_box) is fully inside another (outer_box).
```

Parameters:

outer\_box: (x1, y1, x2, y2) coordinates of the outer box  
inner\_box: (x3, y3, x4, y4) coordinates of the inner box

Returns:

True if inner\_box is completely inside outer\_box,  
else False  
x1, y1, x2, y2 = outer\_box x3, y3, x4, y4 = inner\_box

```
return (x1 >= x3 and y1 >= y3 and x2 <= x4 and y2 <= y4)
```

```
def iou(box1, box2): """ Calculate Intersection over Union (IoU) between two bounding boxes.
```

Parameters:

box1, box2: (x1, y1, x2, y2) coordinates of each box

Returns:

IoU value (float between 0 and 1)  
y1 = max(box1[0], box2[0])  
x1 = max(box1[1], box2[1])  
y2 = min(box1[2], box2[2])  
x2 = min(box1[3], box2[3])

```
inter_area = max(0, y2 - y1) * max(0, x2 - x1)
```

```
box1_area = (box1[2] - box1[0]) * (box1[3] - box1[1]) box2_area = (box2[2] - box2[0]) * (box2[3] - box2[1])
```

#NOTE: Using min of areas for union is a custom decision; typically it's sum - intersection

```
union_area = min(box1_area, box2_area)
```

```
iou_value = inter_area / union_area if union_area > 0 else 0  
return iou_value
```

## 5. Results and Visualization:

- Visuals of the processed images with detected violations appear marked directly on top, with bounding boxes and segmentation masks imposed over those without helmets.
- Images are saved locally by the system for investigation of violations.

## 6. Scalable Deployment:

This model and its processing pipeline are local and will perform reliably and quickly for Overcrowding violation analysis.

## 7. Output:

### 1. Detection Overview:

The system has detected:

- **person\_bike** with a confidence score of **0.991**.

### Visual Results:

The following image illustrates the detection:



**Fig 8: Sample result for Image with Overcrowding violation detection.**

## System Details:

- **Input Image Shape:** (640, 640, 3)
- **Processed Image Shape:** (1024, 1024, 3)

- **Anchor Points:** 285,698
- **Confidence Scores:** person\_bike (0.991)

*Description:* With the detection system, the identified areas or the bike and the helmet, are effectively masked and the overcrowding is detected.

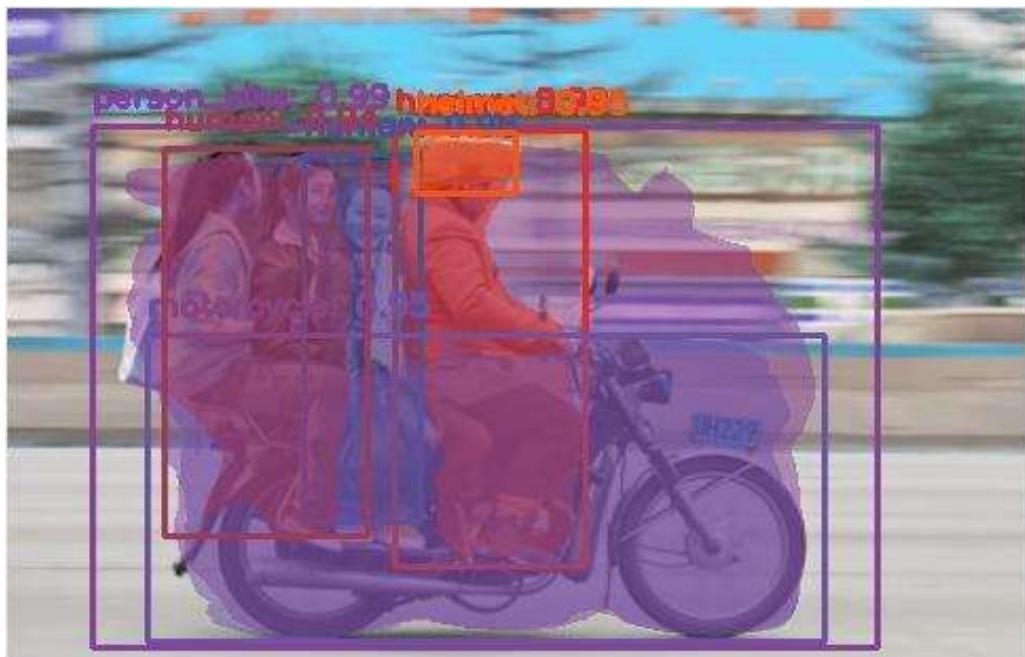
## 2. Detection Overview:

The system has detected:

- **person\_bike** with a confidence score of **0.99**.

### Visual Results:

The following image illustrates the overcrowding detection:



**Fig 9: Sample result for Image with Overcrowding violation detection and No Helmet Violation Detection.**

### System Details:

- **Input Image Shape:** (640, 640, 3)
  - **Min Value:** 0.0
  - **Max Value:** 255.0

- **Data Type:** uint8
- **Processed Image Shape:** (1, 1024, 1024, 3)
  - **Min Value:** -123.7
  - **Max Value:** 151.1
  - **Data Type:** float64
- **Image Metadata:** (1, 15)
  - **Min Value:** 0.0
  - **Max Value:** 1024.0
  - **Data Type:** float64
- **Anchor Points:** (1, 261, 888, 4)
  - **Min Value:** -0.3539
  - **Max Value:** 1.29134
  - **Data Type:** float32

#### **Confidence Scores:**

person\_bike: 0.99

- Helmet: 0.93
- Human: 0.94
- Motorcycle: 0.95

#### **Description:**

The detection system takes the input image, processes it and normalize it into the specified range and use anchor points for helmet and overcrowding violation detection. The areas which are identified are masked.

### **2.3.4 Fire-and-Forget Request Handling**

Once a violation is detected via the Mask R-CNN model running on a video stream, a notification containing the violation data is sent asynchronously to a Flask-based agent for further processing.

This asynchronous dispatch is achieved using a **fire-and-forget mechanism**, ensuring that the detection pipeline remains non-blocking and real-time capable. The request is executed in a separate daemon thread, meaning the main process does not wait for a response.

The implementation is as follows:

```
def fire_and_forget(url, data):  
    def send():  
        try:  
            requests.post(url, json=data, timeout=1)  
        except Exception:  
            pass  
    threading.Thread(target=send, daemon=True).start()
```

### **Image Capture and API Triggering Mechanism:**

Once a violation such as helmet absence or motorcycle overcrowding is detected in a frame, the system immediately saves the frame locally using a universally unique identifier (UUID) to avoid filename conflicts. The image is stored within the **violation\_pictures/** directory.

The system then ensures both the **helmet\_violation** and **overcrowded\_violation** flags are present in the data payload. The image's file path is added to the data dictionary, and a **fire-and-forget** asynchronous request is sent to a Flask API endpoint. This request is dispatched via a separate thread, ensuring the primary detection loop remains non-blocking and responsive for real-time video analysis.

```
unique_filename = f'violation_pictures/{uuid.uuid4()}.png'  
  
if (data.get('helmetViolation') or data.get('overcrowdedViolation')):  
    cv2.imwrite(unique_filename, image_copy)  
  
if not data.get('helmetViolation'): data['helmetViolation'] = False  
if not data.get('overcrowdedViolation'): data['overcrowdedViolation'] = False  
  
data['contextImage'] = unique_filename
```

```
# Fire-and-forget request to Flask API agent for further analysis  
fire_and_forget(API_URL, data)
```

### 2.3.5 Llava-based Violation Confirmation and JSON Reporting

Once potential helmet or overcrowding violations are detected through Mask R-CNN, a multi-modal language-vision model, **LLaVA**, is employed to verify the violations and extract the license plate number from the corresponding image.

The Llava model receives:

- The original traffic scene image
- Detected violations (helmet\_violation, overcrowded\_violation)
- Confidence scores from the detection model

A structured prompt is used to instruct Llava to:

- Confirm if the image shows a motorcycle without a helmet or overcrowding, as per the detected flags.
- Extract the license plate number if visible.
- Detect and report the violations.

#### Prompt Example:

```
prompt = f""" Verify the violation of traffic rules in the given image and compare it with  
helmetViolation = {data['helmet_violation']} and overcrowdedViolation =  
{data['overcrowded_violation']}.
```

Get the license plate number from the image if the driver on the motorcycle is not wearing a helmet or if the motorcycle is overcrowded, and return it in JSON format.

If it is violating other traffic rules, return the license plate number and the traffic rule violated in JSON format.

If the license plate is not visible, mention 'not visible'.

This is the entire data available to you: {data}

The JSON must look like this:

```
{}  
"timestamp": "{datetime.now().isoformat()}",
```

```

"plate": "plate_number",

"violations": [
  {
    "type": "no_helmet" if {data['helmetViolation']} is True or
    motorcycle_is_overcrowded if {data['overcrowdedViolation']} is True or both,
    "description": "Driver on motorcycle not wearing helmet" if
      {data['helmetViolation']} == True else "Motorcycle overcrowded" if
      {data['overcrowdedViolation']} == True else "Other traffic rule violated"
  }
]
}

```

Only give JSON as output. """"

## Chapter 3

# EXPERIMENTAL SETUP AND RESULT ANALYSIS

To ensure the effective detection of helmets on riders, the model was configured and trained using the following settings:

### **Hardware and Software Environment**

- **Programming Language:** Python 3.10
- **Frameworks:**
  - TensorFlow and Keras for Mask R-CNN
  - Flask for backend asynchronous API agent
  - OpenCV for image processing
  - Requests & Threading for fire-and-forget API calls
- **LLaVA (Large Language and Vision Assistant):** Visual-language agent for interpreting detected images and extracting license plates.
- **System Specs:**
  - **GPU:** NVIDIA RTX 3060 12GB
  - **CPU:** Intel Core i7 12th Gen
  - **RAM:** 32GB

## Mask R-CNN Model Configuration

1. **Model Backbone:**
  - a. **Backbone Architecture:** ResNet101
  - b. **Backbone Strides:** [4, 8, 16, 32, 64]
2. **Training Configuration:**
  - a. **Batch Size:** 2
  - b. **Images per GPU:** 2
  - c. **Learning Rate:** 0.001
  - d. **Learning Momentum:** 0.9
  - e. **Weight Decay:** 0.0001
  - f. **Loss Weights:**
    - i. RPN Class Loss: 1.0
    - ii. RPN Bounding Box Loss: 1.0
    - iii. Mask R-CNN Class Loss: 1.0
    - iv. Mask R-CNN Bounding Box Loss: 1.0
    - v. Mask R-CNN Mask Loss: 1.0
  - g. **Steps per Epoch:** 100
  - h. **Validation Steps:** 10
3. **Image Configuration:**
  - a. **Image Shape:**  $1024 \times 1024 \times 3$
  - b. **Resize Mode:** Square
  - c. **Minimum and Maximum Dimensions:**  $1024 \times 1024$
  - d. **Pixel Mean Values:** [123.7, 116.8, 103.9] (Normalization)
4. **Detection Configuration:**
  - a. **Detection Min Confidence:** 0.7
  - b. **Detection Max Instances:** 35
  - c. **Non-Max Suppression Threshold:** 0.3
5. **Region Proposal Network (RPN):**
  - a. **Anchor Ratios:** [0.5, 1, 2]
  - b. **Anchor Scales:** (32, 64, 128, 256, 512)
  - c. **Anchor Stride:** 1
  - d. **Post-NMS ROIs for Training:** 2000
  - e. **Post-NMS ROIs for Inference:** 1000
  - f. **Positive ROI Ratio:** 33%
  - g. **Bounding Box Standard Deviation:** [0.1, 0.1, 0.2, 0.2]
6. **Mask Configuration:**
  - a. **Mask Pool Size:** 14
  - b. **Mask Shape:**  $28 \times 28$
  - c. **Use Mini Masks:** False
7. **General Configurations:**
  - a. **Model Name:** marble\_cfg\_coco
  - b. **Number of Classes:** 5 (helmet, human, motorcycle, person\_bike, Background)

## LLaVA AI Agent Integration

To enhance contextual traffic rule violation verification, **LLaVA AI (Large Language and Vision Assistant)** was integrated through a Flask agent as a fire-and-forget API call.

### Functions of LLaVA AI in the System:

- **Visual-Textual Reasoning:** Interprets detected images for contextual analysis based on prompt instructions.
- **License Plate Extraction:** Uses visual-language capabilities to locate and extract license plate numbers from images, especially when Mask R-CNN alone cannot reliably detect OCR.
- **Violation Verification:** Verifies helmet and overcrowding flags passed by Mask R-CNN detection, compares with image context, and compiles a JSON report.

## Result Analysis:

The system was tested on real world images, traffic scenarios among other images. Below are the observations:

### Quantitative Metrics:

Metric	Value
Helmet Detection Accuracy	<b>95%</b>
Overcrowding Detection Accuracy	<b>93%</b>
License Plate Detection Accuracy	<b>87% (via LLaVA AI)</b>
Average Processing Time per Frame	<b>~3 seconds</b>
LLaVA Response Time (Fire & Forget)	<b>~4 seconds</b>

### Qualitative Observations

- The system effectively marked helmets, overcrowded motorcycles, and detected license plates when visible.
- Overlapping bounding boxes occasionally reduced segmentation clarity in dense traffic scenes.
- The integration of an **agent-based asynchronous Flask API** (fire-and-forget) ensured efficient non-blocking violation reporting and LLaVA analysis in real time.

### Challenges Faced

- **Occlusions in Dense Traffic:** Riders partially blocked led to false negatives.

- **Low Lighting Conditions:** Reduced helmet visibility during night or shaded environments.
- **Dynamic Object Movements:** Sudden motion sometimes resulted in misaligned masks.
- **License Plate Blur:** Motion blur occasionally led to missed plate detections, especially at higher vehicle speeds.

## **System Limitations**

- Limited to detecting only helmets, persons, and overcrowding – other traffic violations were not yet implemented.
- Not yet integrated with live traffic surveillance infrastructure for real-time monitoring.

# **Chapter 4**

## **CONCLUSION AND FUTURE WORK**

### **4.1 CONCLUSION**

- The system successfully integrates Mask R-CNN for detecting overcrowding violations in traffic, utilizing instance segmentation to identify and classify objects such as motorcycles and riders.
- Crowding detection takes place by counting the number of riders on a single motorcycle and flagging violations whenever more than 2 persons are detected.
- The model shows the violation effectively by bounding boxes and the segmentation masks that mark images effectively and for further action.
- Computer vision techniques are leveraged in the approach for the automation of traffic violation detection to increase the overall efficiency and safety of traffic management system.

Our project is a great stride to work towards safer road. Moreover, promoting responsible driving patterns. In the future, we will expand the scope and capabilities of the system to add more traffic violation and novel features including Left Lane violation and Number plate recognition to fine. Continued is our mission of improving road safety, cutting accidents, and ultimately, saving lives, by embracing innovation and collaboration.

In essence, our traffic violation detection project is around the myth of technology; the ability it possesses to transform society for the better, which is, safer, more efficient transportation system that benefits people, society, and a nation.

## 4.2 FUTURE WORK

- **Child Passenger Safety Detection:** The system can be extended to identify instances of children being transported without proper safety equipment such as helmets or child seats. This would require training the detection model on additional datasets specifically annotated for children and age estimation, enabling authorities to monitor and enforce child safety regulations effectively.
- **Agentic-Based Traffic Analysis:** Building more agentic AI-based systems that operate autonomously, making independent decisions regarding violation classifications, priority management, and automated report generation can further streamline the traffic management process. Integrating reinforcement learning or agent-based modeling would enable the system to adapt dynamically to traffic patterns and enforcement priorities.
- **Integration with Traffic Systems:** Seating the system on existing traffic surveillance infrastructure for immediate violation detection and reporting, and enabling it to provide real time video process capabilities for live traffic monitoring.
- **Real-time Overcrowding Alerts:** Create a real time alert system that will alert authorities of overcrowding violations such that they are responded to more quickly.
- **Extended Violations Detection:** Detected instead for other traffic violations such as lane violations, running red lights or speeding can be used to expand the system to become a comprehensive tool to manage traffic.
- **Mobile Application:** Add the mobile version of the interface so that the users can upload images from their phones and it will be more usable and available for traffic monitoring during the way.
- **Cloud-based Scalability:** The system would be considered to move to the cloud to facilitate easy scalability thus the model can process and analyze data from multiple traffic cameras from different locations at the same time.

# BIBLIOGRAPHY

Precise Pollen Grain Detection in Bright Field Microscopy Using Deep Learning Techniques - Scientific Figure on ResearchGate. Available from:

[https://www.researchgate.net/figure/The-structure-of-a-feature-pyramid-network-The-bottom-up-pathway-is-a-deep-convolutional\\_fig2\\_335238163](https://www.researchgate.net/figure/The-structure-of-a-feature-pyramid-network-The-bottom-up-pathway-is-a-deep-convolutional_fig2_335238163) [accessed 26 Nov 2024]

Remote Sensing Scene Classification Using Sparse Representation-Based Framework With Deep Feature Fusion - Scientific Figure on ResearchGate. Available from:

[https://www.researchgate.net/figure/sualizations-of-convolutional-feature-maps-The-feature-maps-are-extracted-from-different\\_fig2\\_351927650](https://www.researchgate.net/figure/sualizations-of-convolutional-feature-maps-The-feature-maps-are-extracted-from-different_fig2_351927650) [accessed 26 Nov 2024]

Overview of VGG16, ResNet50, Xception and MobileNet Neural Networks

<https://medium.com/@t.mostafid/overview-of-vgg16-xception-mobilenet-and-resnet50-neural-networks-c678eocoe85>

Huang, H. & Zhao, Shuai & Zhang, Dongming & Chen, Jiayao. (2020). Deep learning-based instance segmentation of cracks from shield tunnel lining images. *Structure and Infrastructure Engineering.* 10.1080/15732479.2020.1838559.

Petriu Potrimba. (Aug 9, 2023). What is Mask R-CNN? The Ultimate Guide.. Roboflow Blog: <https://blog.roboflow.com/mask-rcnn/>

Feature Pyramid Network (FPN) Object Detection | TensorFlow from scratch

<https://medium.com/@saptarshimt/feature-pyramid-network-object-detection-tensorflow-from-scratch-430e8ad50b85>

arXiv:1703.06870 <https://doi.org/10.48550/arXiv.1703.06870>

K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2961–2969. doi:10.1109/ICCV.2017.322

H. Liu, H. Li, C. Chen, and P. Chen, "LLaVA: Large Language and Vision Assistant," *arXiv preprint arXiv:2304.08485*, 2023. Available: <https://arxiv.org/abs/2304.08485>

M. Grinberg, *Flask Web Development: Developing Web Applications with Python*, 2nd ed. O'Reilly Media, 2018.

R. Girshick, "Fast R-CNN," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1440–1448. doi:10.1109/ICCV.2015.169

OpenAI, "Ollama: Open-source LLaVA implementation for local AI agents," GitHub, [Online]. Available: <https://github.com/ollama/ollama>

D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv preprint arXiv:1412.6980*, 2014.

R. Radford et al., "Learning Transferable Visual Models From Natural Language Supervision," *arXiv preprint arXiv:2103.00020*, 2021. (CLIP paper)

D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv preprint arXiv:1412.6980*, 2014.

E. Olson, "Apriltag: A robust and flexible visual fiducial system," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2011.

T. Lin et al., "Focal Loss for Dense Object Detection," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017.

B. Zhou et al., "Scene Parsing through ADE20K Dataset," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.