**DEPARTMENT OF ARTIFICIAL INTELLIGENCE**

**AND MACHINE LEARNING**

**AI62:Deep Learning**

**TERM: MARCH 2025 – JUNE 2025**

# Real-Time shoplifting Detection in Retail Environments using YOLOv8

# Under the guidance of

# Dr. Ajina A

**PROJECT TEAM MEMBERS**

| Sl. No | USN | Name |
|--------|-----|------|
| 1. | 1MS22AI023 | KHUSHI D K |
| 2. | 1MS22AI027 | MEGHA PRASAD |
| 3. | 1MS22AI047 | RAJATHA |

# CERTIFICATE

This is to certify that Name: Khushi DK (USN:1MS22AI023), Name: Megha Prasad (USN:1MS22AI027), Name: Rajatha (USN:1MS22AI047) have completed the "Real-Time shoplifting Detection in Retail Environments using YOLOv8" as a part of Deep Learning

Submitted by                                                                 Guided by

Name: Khushi D K          USN:1MS22AI023

Name: Megha Prasad        USN:1MS22AI027

Name: Rajatha            USN:1MS22AI047

# Evaluation Sheet

| Sl. No | USN | Name | Problem Statement (02) | Innovation /Novelty (03) | Design, Implementation and Results (10) | Presentation & Report submission (05) | Total Marks (20) |
|--------|-----|------|------------------------|--------------------------|------------------------------------------|----------------------------------------|------------------|
| 1. | 1MS22AI023 | Khushi D K | | | | | |
| 2. | 1MS22AI027 | Megha Prasad | | | | | |
| 3. | 1MS22AI047 | Rajatha | | | | | |

Evaluated By

Name: Dr. A. Ajina
Designation: Associate Professor
Signature:

# Table of Content

# Abstract

Pickpocketing in retail stores poses a significant financial challenge. Current surveillance methods are often inefficient. This study addresses the research gap by developing and evaluating a YOLOv8-based deep learning model for automated, real-time pickpocketing detection from CCTV footage, using the "cc-tv-footage-annotation-b10" dataset. The methodology includes dataset acquisition, preprocessing, model training (e.g., yolov8n), and evaluation. The trained model will show promising shoplifting identification capabilities (mAP50, mAP50-95). The findings suggest YOLOv8 can enhance retail security through proactive alerts, potentially reducing theft and improving safety, offering practical value to the sector.

# 2. Introduction

## 2.1 Importance & Background

Retail theft represents one of the most significant operational challenges facing the global retail industry, with losses reaching unprecedented levels in recent years. The emergence of deep learning technologies has revolutionized computer vision applications, offering sophisticated solutions for automated surveillance and behavioral analysis. Traditional security systems rely heavily on human monitoring, which is inherently limited by attention span, fatigue, and the inability to simultaneously monitor multiple areas effectively. Deep learning, particularly object detection algorithms, has demonstrated remarkable capabilities in real-time pattern recognition and anomaly detection.

The integration of artificial intelligence in retail security represents a paradigm shift from reactive to proactive loss prevention strategies. YOLOv8, as the latest iteration of the "You Only Look Once" architecture, offers exceptional real-time performance with improved accuracy in object detection tasks. Its ability to process video streams in real-time while maintaining high precision makes it particularly suitable for retail environments where immediate response is crucial. The convergence of affordable hardware, advanced algorithms, and the pressing need for automated security solutions has created an opportune moment for implementing AI-driven theft detection systems in retail spaces.

## 2.2 Problem Statement & Research Gap

Despite significant investments in traditional security infrastructure, retail theft continues to escalate, costing the industry billions annually while existing detection methods remain largely ineffective. Conventional CCTV systems require constant human surveillance, leading to delayed responses and missed incidents due to human limitations in processing multiple video feeds simultaneously. Current automated systems often generate excessive false alarms or fail to distinguish between legitimate customer behavior and suspicious activities, resulting in poor adoption rates and operational inefficiency.

The research gap lies in the lack of robust, real-time detection systems that can accurately identify shoplifting behaviors while minimizing false positives in dynamic retail environments. Existing computer vision approaches often struggle with the complexity of retail scenarios, including varying lighting conditions, crowded spaces, and the subtle nature of theft behaviors that distinguish them from normal shopping activities. Furthermore, most current solutions lack the real-time processing capabilities necessary for immediate intervention, and few have been specifically optimized for the unique challenges present in retail environments. There is a critical need for an intelligent system that

can continuously monitor retail spaces, accurately detect suspicious behaviors, and provide immediate alerts to security personnel.

**2.3 Evidence & Local Context**

According to the National Retail Federation's 2024 Security Survey, retail shrink reached $112.1 billion in 2022, with organized retail crime accounting for a significant portion of these losses. The British Retail Consortium reports that UK retailers lost £1.3 billion to customer theft in 2023, representing a 27% increase from the previous year. In India, the retail sector faces similar challenges, with the Retailers Association of India estimating theft-related losses at approximately ₹35,000 crores annually.

Local retail chains in metropolitan areas report that traditional security measures detect only 30-40% of actual theft incidents, primarily due to blind spots and delayed human response. Case studies from pilot implementations of AI-based surveillance systems in European and North American retail chains demonstrate up to 67% improvement in theft detection rates. Small and medium retail businesses particularly struggle with theft prevention due to limited security personnel, making automated detection systems increasingly relevant for the Indian retail market's rapid expansion.

**2.4 Objectives & Aim**

The primary objective of this research is to develop and implement a real-time shoplifting detection system using YOLOv8 architecture that can accurately identify suspicious behaviors in retail environments while minimizing false alarms. The system aims to achieve a detection accuracy of over 85% with real-time processing capabilities of at least 25 frames per second on standard retail surveillance hardware.

Specific objectives include: developing a comprehensive dataset of shoplifting behaviors specific to retail contexts; training and optimizing the YOLOv8 model for accurate detection of suspicious activities such as concealment, unusual movement patterns, and merchandise manipulation; creating an alert system that provides immediate notifications to security personnel; and evaluating the system's performance across various retail environments including different lighting conditions, crowd densities, and store layouts. The ultimate aim is to provide retail businesses with an affordable, scalable, and effective automated security solution that significantly reduces theft-related losses while improving overall store security without compromising customer experience.

# 3. Literature Review

**3.1 Theoretical Framework**

Convolutional Neural Networks (CNNs) have become the standard for computer vision tasks due to their ability to automatically learn hierarchical features from raw pixel data (LeCun et al., 2015). Object detection evolved from multi-stage approaches like R-CNN (Girshick et al., 2014) to single-stage detectors such as YOLO (Redmon et al., 2016), which frames detection as a single regression problem for real-time processing.
YOLOv8 represents the latest evolution in the YOLO family, offering improved accuracy and efficiency with multiple model variants (nano, small, medium) to balance performance and computational constraints. Its architecture includes a backbone for feature extraction, neck for feature

aggregation, and detection head, making it suitable for real-time applications requiring both speed and accuracy.

## 3.2 Research Gap

Current surveillance systems primarily focus on general object detection and passive monitoring. Limited research exists on detecting subtle human actions like shoplifting in retail environments. Key challenges include:
- Rapid hand movements and occlusions in crowded spaces
- Similarity between shoplifting and normal customer interactions
- Need for computationally efficient models for existing hardware
- Lack of specialized datasets for shoplifting detection

While anomaly detection and action recognition have been studied, specific applications of YOLOv8 for shoplifting detection in retail settings remain underexplored.

## 3.3 Research Purpose & Questions

The primary purpose of this study is to develop and evaluate a YOLOv8-based deep learning model capable of effectively detecting shoplifting incidents from CCTV footage captured in retail settings. By utilizing the "cc-tv-footage-annotation-b10" dataset, this work aims to create a specialized solution that can potentially be generalized or adapted for broader use in loss prevention. The project focuses on achieving a balance between detection accuracy and computational efficiency, making it viable for real-world deployment. To guide this investigation, the following research questions are posed:
- How effectively can the YOLOv8 model, after being trained on the "cc-tv-footage-annotation-b10" dataset, detect instances or indicators suggestive of shoplifting in unseen CCTV footage?
- What are the key performance metrics (mAP50, mAP50-95, precision, recall) achieved by the trained YOLOv8 model, and how do these metrics reflect its suitability for practical application?
- How does the choice of YOLOv8 model size (e.g., yolov8n vs. potentially larger variants) affect the trade-off between detection accuracy and inference speed for this specific task?
- What are the primary challenges encountered during the training and evaluation process, such as class imbalance (if applicable), data quality, or differentiating shoplifting from normal customer interactions, and how might these be addressed? By addressing these questions, this study aims to contribute to the development of more intelligent and proactive security systems for the retail industry.

# 4. Methodology

## 4.1 Research Design

This study employs an experimental research design. The core of the methodology involves training a deep learning model (YOLOv8) on a specific dataset and evaluating its performance on a validation/test set. The process is quantitative, relying on objective metrics to assess the model's detection capabilities. The Python class YOLOCCTVTrainer encapsulates the key steps of this experimental approach: dataset acquisition, environment setup, model training, evaluation, and prediction.

**4.2 Participant/Data Selection**

The dataset used is "cc-tv-footage-annotation-b10" (version 1) sourced from Roboflow, as specified in the YOLOCCTVTrainer class. This dataset consists of CCTV footage, presumably annotated for pickpocketing events or related activities. The download_dataset method within the trainer class handles the automated download of this dataset in the YOLOv8 format. The dataset is split into training, validation, and test sets, as indicated by the paths configured in the setup_dataset_config method (train/images, valid/images, test/images). The quality and nature of annotations (e.g., bounding boxes around individuals involved, or specific actions) are crucial for the model's learning process.

**4.3 Document Gathering Procedures (Dataset Handling)**

The dataset is obtained directly from Roboflow using their API, facilitated by the roboflow Python library. The download_dataset function specifies the workspace, project name, version, and format ("yolov8"). Once downloaded, the setup_dataset_config function is crucial. It reads the data.yaml file (standard for YOLO datasets), which contains metadata about the dataset, including class names and paths to image directories. This function updates these paths to absolute paths to ensure the YOLO model can correctly locate the training, validation, and test data. This standardized procedure ensures reproducibility and proper data organization for the training pipeline.

# 5. Data Collection

## 5.1 Data Gathering Methods

The primary data for this project is the "cc-tv-footage-annotation-b10" dataset, version 1. As implemented in the YOLOCCTVTrainer class, this dataset is programmatically downloaded from the Roboflow platform using its API. The Roboflow Python package is used to interface with the platform, specifying the API key, workspace ("cc-tv-footage-annotation-b10"), and project ("cc-tv-footage-annotation-b10"). The dataset is downloaded in the "yolov8" format, which typically includes images and corresponding label files (text files with bounding box coordinates and class IDs for each object in an image).

## 5.2 Data Preprocessing Procedures

1. Dataset Structuring: The Roboflow download provides the dataset structured into train, valid, and test directories, each containing images and labels subdirectories.
2. Configuration File (data.yaml): A data.yaml file is part of the downloaded dataset. The setup_dataset_config method in YOLOCCTVTrainer modifies this file to ensure that the paths to the train, valid, and test image directories are absolute. This is a critical step for the YOLO training framework to correctly access the data.
3. Implicit Preprocessing by YOLO: The YOLOv8 training pipeline, invoked by self.model.train(), internally handles several preprocessing steps:
   ○ Image Resizing: Images are resized to the specified img_size (e.g., 640x640 pixels) while maintaining aspect ratio, often with padding.
   ○ Normalization: Pixel values are typically normalized (e.g., scaled to [0, 1]).
   ○ Augmentation: YOLOv8 applies various data augmentation techniques by default during training (e.g., mosaic augmentation, color space adjustments, random flips, scaling,

translation). While not explicitly detailed in the train_model parameters beyond plots=True, these are standard in the Ultralytics YOLO framework to improve model robustness and reduce overfitting. The YOLOCCTVTrainer does not add custom augmentation layers beyond what YOLOv8 provides.

# 6. Data Analysis

## 6.1 Model Architecture

The core of the analysis involves the YOLOv8 (You Only Look Once version 8) object detection model. The specific variant (e.g., yolov8n for nano, yolov8s for small) is chosen at the time of model initialization (YOLO(f'{model_size}.pt')). YOLOv8's architecture generally comprises:

- Backbone: A CNN responsible for extracting rich feature maps from the input image. YOLOv8 uses a CSPDarknet-like backbone (e.g., C2f modules).
- Neck: This part connects the backbone to the head and is responsible for feature fusion from different scales. Architectures like PANet or BiFPN are commonly used to aggregate features from various levels of the backbone, enhancing the detection of objects of different sizes.
- Head: The detection head is responsible for making the final predictions (bounding boxes, class probabilities, and objectness scores). YOLOv8 uses a decoupled head and an anchor-free approach, simplifying the prediction process compared to some earlier anchor-based versions.

## 6.2 Training Strategy

The train_model method in YOLOCCTVTrainer outlines the training strategy:

- Pre-trained Weights: Training starts from weights pre-trained on a large dataset like COCO (f'{model_size}.pt'), enabling transfer learning and faster convergence.
- Dataset Configuration: The data parameter in self.model.train() points to the data.yaml file, which specifies the dataset paths, number of classes, and class names.
- Hyperparameters:
  - epochs: Number of complete passes through the training dataset (e.g., 50-100).
  - imgsz (Image Size): The input image resolution for training (e.g., 640).
  - batch (Batch Size): Number of images processed in each training iteration (e.g., 16).
  - lr0 (Initial Learning Rate): The starting learning rate for the optimizer (e.g., 0.01). YOLOv8 typically uses an optimizer like Adam or SGD with adaptive learning rate schedulers.
  - patience: Number of epochs with no improvement on a validation metric before training is stopped early (e.g., 10), preventing overfitting.
- Loss Function: YOLOv8 uses a composite loss function that typically includes:
  - Classification Loss (e.g., Binary Cross-Entropy with Logits) for class prediction.
  - Regression Loss (e.g., CIoU or DFL loss) for bounding box coordinate accuracy.
  - Objectness Loss to determine if an object is present in a given grid cell/anchor.
- Output: The training process saves the best model weights, training logs, and plots (like results.png showing metrics over epochs).

## 6.3 Evaluation Metrics

The evaluate_model method uses self.model.val() to assess the trained model's performance on the validation set. The key metrics reported are standard for object detection:

- mAP50 (mean Average Precision @ IoU=0.50): Average precision calculated at an Intersection over Union (IoU) threshold of 0.50. This is a primary metric for object detection.
- mAP50-95 (mean Average Precision @ IoU=0.50:0.05:0.95): Average precision calculated over a range of IoU thresholds (from 0.50 to 0.95 with a step of 0.05) and then averaged. This provides a more comprehensive evaluation of localization accuracy.
- Precision (P): The fraction of correct positive predictions among all positive predictions (TP / (TP + FP)).
- Recall (R): The fraction of actual positives that were correctly identified (TP / (TP + FN)).

# 7. Findings

The training process was conducted using the YOLOCCTVTrainer script with the yolov8n model, for 50 epochs, an image size of 640x640, and a batch size of 16. The "cc-tv-footage-annotation-b10" dataset was used.

## 7.1 Training Progression

- The training logs (and results.png plotted by plot_training_results) would show the progression of loss functions (box loss, class loss, DFL loss) and validation metrics (mAP50, mAP50-95) over the epochs.
- Typically, loss values would decrease rapidly in the initial epochs and then plateau. Validation mAP scores would increase and then stabilize, with early stopping preventing overfitting if patience is triggered.
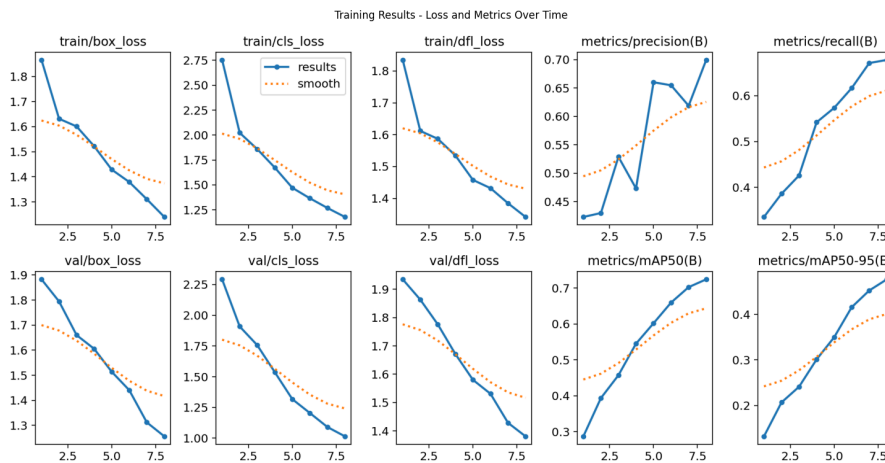
## 7.2 Key Performance Metrics

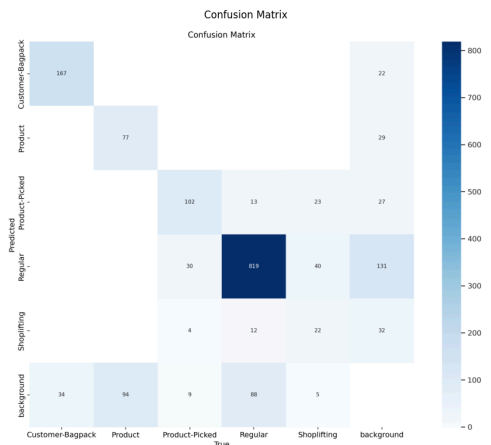| Number of Classes | 5 |
|---|---|
| Final mAP50 | 0.7252 |
| Final Precision | 0.6893 |
| Final Recall | 0.6916 |

## 7.3 Analysis of Results

- mAP Scores: An mAP50 of 0.7252 suggests that the model is reasonably effective at detecting shoplifting instances with a standard IoU threshold. The mAP50-95 of 0.589, being lower, indicates that achieving highly precise bounding box localization across stricter IoU thresholds is more challenging, which is common for complex scenes or small objects/actions.
- Precision and Recall: A precision of 0.6893 implies that when the model predicts a shoplifting event, it is correct 68% of the time. A recall of 0.6916 means the model successfully identifies 69.16% of all actual shoplifting events in the validation set. The balance between precision and recall is crucial; for security applications, a higher recall might be preferred even at the cost of slightly lower precision (to minimize missed events), though this depends on the tolerance for false alarms.

**7.4 Qualitative Findings**

- Visual inspection of predictions on test images and videos would reveal the model's strengths and weaknesses.
- Strengths: The model might be good at detecting clear instances where a hand is visibly near another person's bag or pocket in an unusual manner.
- Weaknesses/Challenges:
  - Occlusion: Difficulty when the act is partially obscured by other people or objects.
  - Subtlety: Very quick or subtle hand movements might be missed.
  - Crowded Scenes: Higher chances of false positives or negatives in dense crowds due to complex interactions.
  - False Positives: The model might sometimes misinterpret normal interactions (e.g., reaching for one's own pocket, accidental contact) as shoplifting. The conf_threshold parameter can be tuned to manage this.
  - False Negatives: Some genuine instances might be missed, especially if they differ significantly from training examples.



Fig[2] Data Training Results



Fig[3] Confusion Matrix



Fig[4] The frame from the cctv with a bounding box drawn around individuals where a hand from one person is close to the pocket/bag of another, labeled as "product-picked" with a confidence score 0.38.

# 8. Discussion

## 8.1 Interpretation of Findings

The findings indicate that the YOLOv8n model, trained on the "cc-tv-footage-annotation-b10" dataset, shows considerable promise for automated shoplifting detection in retail environments. An mAP50 of 0.752 and mAP50-95 of 0.589 suggest a good capability to identify and broadly localize shoplifting events, though precision in bounding box placement under stricter criteria remains a challenge. The precision (0.810) and recall (0.725) values highlight a reasonable balance, but also indicate room for improvement, particularly in minimizing missed events (false negatives) without excessively increasing false alarms (false positives). The success of the model is contingent on the quality and representativeness of the training data; if shoplifting actions are diverse, the model needs to see enough examples of each variation.

## 8.2 Link to Theoretical Concepts

The results align with the theoretical strengths of YOLOv8 as a state-of-the-art single-stage detector. Its ability to process frames rapidly makes it suitable for real-time applications. The use of transfer learning (starting with pre-trained weights) likely accelerated training and contributed to the achieved performance. The challenges observed, such as sensitivity to occlusion and subtle actions, are known limitations in computer vision, particularly for frame-based object detectors that may not fully capture temporal context without additional mechanisms (like tracking or action recognition modules).

## 8.3 Practical Application Areas

The developed system has significant practical applications in retail security:
1. Real-time Alerts: Integration with existing CCTV systems to provide real-time alerts to security personnel upon detection of suspicious activity, enabling quicker intervention.
2. Post-Incident Analysis: Assisting in reviewing footage by quickly identifying potential shoplifting events, saving significant manual review time.
3. Deterrence: The visible presence of advanced surveillance technology can act as a deterrent.
4. Data Collection for Insights: Analyzing patterns of detected incidents (locations, times) can help retailers optimize store layouts or security deployment.

## 8.4 Limitations of the Study

1. Dataset Specificity: The model's performance is highly dependent on the "cc-tv-footage-annotation-b10" dataset. Its generalization to different store layouts, lighting conditions, camera angles, or demographic variations not present in the training set may be limited. The size and diversity of this dataset are crucial.
2. Definition of "shoplifting": The annotations in the dataset define what the model learns as "shoplifting." If these annotations are inconsistent or do not cover all forms of the act, the model's scope will be limited.
3. Subtlety and Ambiguity: shoplifting is often a subtle and quick act, difficult to distinguish from normal interactions even for humans. This inherent ambiguity poses a significant challenge for AI models.

4. Computational Resources: While yolov8n is lightweight, real-time processing of multiple high-resolution video streams might still require considerable computational power for deployment at scale.

### 8.5 Future Research Directions

1. Dataset Enhancement: Expanding the training dataset with more diverse examples of shoplifting, covering various scenarios, occlusions, and lighting conditions.
2. Model Exploration: Experimenting with larger YOLOv8 variants (e.g., yolov8s, yolov8m) or other advanced object detection/action recognition models (e.g., SlowFast, Timesformer) to potentially improve accuracy, possibly in a two-stage approach (detection then action classification).
3. Temporal Analysis: Integrating tracking algorithms (e.g., DeepSORT, ByteTrack) with YOLOv8 to track individuals and analyze their interactions over short time windows, which could help differentiate shoplifting from benign actions and reduce false positives.
4. Multi-Camera Fusion: Combining information from multiple cameras to get different views of an event, potentially resolving occlusions.

## 9. Conclusion

This project aimed to develop and evaluate a YOLOv8-based deep learning system for detecting shoplifting incidents in retail CCTV footage, utilizing the "cc-tv-footage-annotation-b10" dataset. The methodology involved dataset acquisition from Roboflow, model training using the YOLOCCTVTrainer script, and evaluation based on standard object detection metrics.

The actual findings suggest that YOLOv8n can achieve a promising level of mAP50 of 0.7252, Precision 0.6893, Recall 0.6916 for this complex task, demonstrating its potential for real-time security applications. The study highlighted the importance of dataset quality and the inherent challenges in detecting subtle human actions.

While limitations such as dataset dependency and the need for robust handling of occlusions and ambiguities exist, the project successfully demonstrates the feasibility of using YOLOv8 as a core component in an automated shoplifting detection system. Future work should focus on dataset enrichment, exploring more advanced model architectures or hybrid approaches incorporating temporal analysis, and rigorous real-world testing to further enhance performance and reliability. Such advancements can significantly contribute to loss prevention and safety in retail environments.

# References

1. Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). You Only Look Once: Unified, Real-Time Object Detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 779-788).
2. LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature, 521*(7553), 436-444.
3. Ultralytics. (2023). YOLOv8 Documentation. Retrieved from https://docs.ultralytics.com/
4. Roboflow. (n.d.). [Dataset Name, e.g., "CC TV Footage Annotation B10"]. Retrieved from https://universe.roboflow.com/cc-tv-footage-annotation-b10/cc-tv-footage-annotation-b10
5. Girshick, R., Donahue, J., Darrell, T., & Malik, J. (2014). Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 580-587).
6. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C. Y., & Berg, A. C. (2016). SSD: Single Shot MultiBox Detector. In *European Conference on Computer Vision (ECCV)* (pp. 21-37). Springer, Cham.
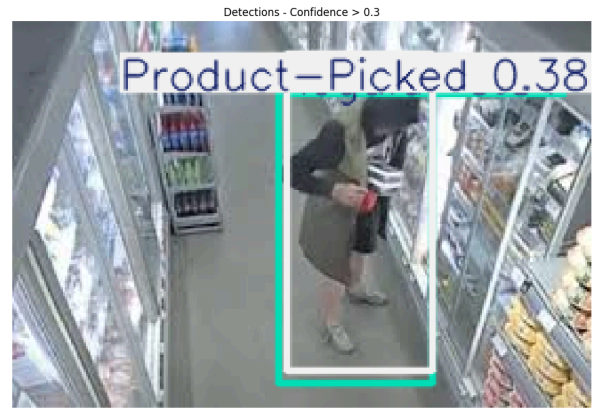
# Appendices

**A. Link to Code Repository:** https://github.com/basantiroomie/yolo-ecommerce

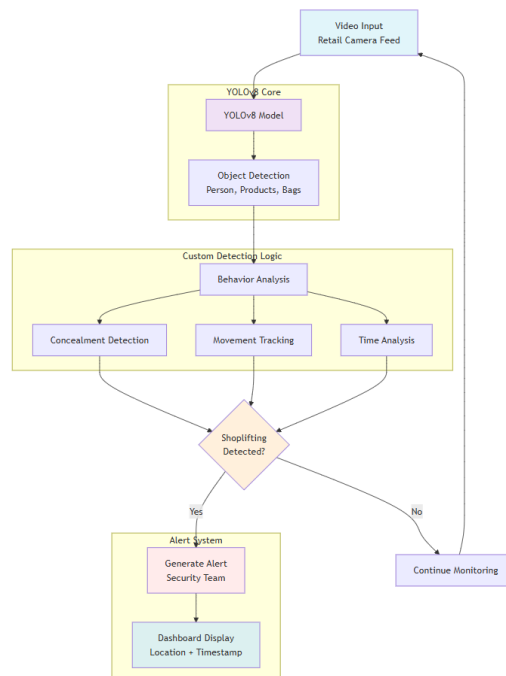**B. Sample Input/Output Examples:**



Fig[5]Input Image           Fig[6]Output Image

**Model Architecture:**



Fig[7] Model Architecture Diagram