

# AI-Based Surveillance System For Abandoned Object Detection Using YOLOv8

Mr. Nithishkumar P<sup>1</sup>, Sakthi S<sup>2</sup>, Praveen S<sup>3</sup> MS. Deepa G<sup>4</sup>

<sup>1</sup>Assist prof, Dept of Artificial Intelligence and Data Science

<sup>2,3,4</sup>Dept of Artificial Intelligence and Data Science

<sup>1,2,3,4</sup> Sri Venkateswara College of Engineering

**Abstract-** *The proliferation of surveillance infrastructure across transportation hubs, commercial complexes, and civic spaces has not been matched by a proportional improvement in the capacity of human operators to effectively monitor multiple video feeds over extended durations. Cognitive limitations inherent to sustained visual monitoring create critical gaps during which security-relevant events, including the placement of unattended items, may go unnoticed. This paper proposes a computationally efficient, learning-based framework that autonomously identifies abandoned personal belongings within live surveillance footage by integrating the YOLOv8 single-stage detector with a spatiotemporal ownership inference mechanism. The system processes individual video frames to simultaneously detect persons and personal items — encompassing backpacks, handbags, laptops, and mobile devices — using COCO-pretrained weights applied in a zero-shot configuration. Temporal continuity is preserved through an IoU matching strategy and ownership attribution is determined by Euclidean proximity analysis. Should the attended condition remain unsatisfied beyond a user-defined temporal threshold, the item is reclassified as abandoned, prompting a visual alert. Experimental validation confirms near-real-time throughput and reliable detection outcomes.*

**Keywords:** Abandoned Object Detection, YOLOv8 Deep Learning, Spatiotemporal Ownership Inference, Intersection-over-Union Tracking, Zero-Shot Deployment, Euclidean Proximity Analysis, Real-Time CCTV Surveillance, Public Safety Automation.

## I. INTRODUCTION

Today, cameras are installed almost everywhere — airports, train stations, shopping malls, schools, and public parks. These cameras produce a large amount of video footage every single day. However, watching all this footage manually is not practical. Security guards cannot monitor so many screens at the same time for long hours without making mistakes. One important type of security threat that often goes unnoticed is an abandoned bag or object left behind by someone in a public place. Such objects can be dangerous,

stolen, or simply lost — and identifying them quickly is critical for public safety.

Currently, most surveillance systems depend entirely on human operators to spot suspicious objects on camera. A person watching multiple screens for several hours will naturally become tired and lose focus over time. To solve this problem, there is a need for an automated system that can continuously watch the video feed and automatically detect when an object has been left behind without human effort.

In recent years, deep learning technology has made it possible to detect objects in video very accurately and very fast. One of the best tools available today is YOLOv8, a deep learning model developed by Ultralytics. YOLOv8 can detect many types of objects — including people, bags, laptops, and phones — in real time from a camera feed. It is fast enough to process video frames as they come in, making it ideal for live surveillance applications.

This paper presents an intelligent surveillance system that uses YOLOv8 to automatically detect abandoned objects in video. The system works in three steps: first it detects all persons and personal belongings in each frame; second it tracks each detected object across multiple frames; third it checks whether any person is standing near the object. If an object stays in one place for more than a set number of seconds and no person is nearby, the system marks it as abandoned and shows an alert on the screen along with a timestamped incident log entry.

## II. LITERATURE SURVEY

Several researchers have worked on the problem of detecting suspicious or abandoned objects in surveillance videos. These efforts can be broadly grouped into four areas: (1) traditional background subtraction methods, (2) deep learning-based object detection, (3) multi-object tracking techniques, and (4) dedicated abandoned object detection systems.

Background Subtraction Methods

The earliest automated surveillance systems used background subtraction, where the system builds a model of the "empty" scene and identifies anything that does not belong. Stauffer and Grimson [3] proposed the Gaussian Mixture Model (GMM), which became the standard background subtraction method for many years. Piccardi [4] reviewed various techniques and noted that all of them perform poorly when the background itself changes — due to flickering lights, shadows, or moving trees — making them unreliable for busy public spaces.

#### Deep Learning-Based Object Detection

The introduction of deep learning transformed object detection completely. Girshick et al. [5] proposed R-CNN, the first method to use deep features for detection. Faster R-CNN [6] improved speed significantly. The YOLO family, introduced by Redmon et al. [7], solved the real-time challenge by processing the entire image in a single pass. The latest version, YOLOv8 [1], uses an anchor-free detection head and improved training methods, making it the most suitable backbone for the proposed surveillance system.

#### Multi-Object Tracking Techniques

Detecting an object in one frame is not enough — a surveillance system also needs to track the same object across multiple frames to measure how long it has been stationary. Bewley et al. [9] proposed SORT, which uses Kalman filtering and the Hungarian algorithm. Wojke et al. [10] extended SORT into DeepSORT with a deep appearance network. The proposed system uses a simpler IoU-based matching approach, which is sufficient for tracking stationary objects while avoiding the overhead of deep re-identification.

#### Abandoned Object Detection Systems

Some researchers have specifically studied abandoned object detection in surveillance videos, typically combining object detection with a timer-based ownership rule. However, many of these systems rely on separate background subtraction and deep learning stages, making them complex to deploy. The proposed system simplifies this by using a single YOLOv8 model to handle both person and object detection simultaneously, removing the need for any separate segmentation stage.

#### Research Gap Addressed by this Work

A review of existing literature shows that most abandoned object detection systems either rely on fragile background subtraction methods or require complex multi-

stage pipelines. Very few systems combine modern single-stage object detection with a lightweight spatial proximity rule in a zero-shot configuration. This paper directly addresses this gap by proposing an end-to-end pipeline that is accurate, fast, and deployable on standard hardware without any fine-tuning.

### III. PROPOSED SYSTEM ARCHITECTURE

The proposed system is designed to automatically detect abandoned objects in real-time surveillance video without requiring any human involvement. The entire pipeline is built around a single deep learning model — YOLOv8 — which handles all object and person detection tasks. The output is then processed by five additional modules that together determine whether a detected object has been abandoned. The system processes video frame by frame, detects all persons and personal belongings, tracks them across frames, and checks whether any person is nearby.

#### Object Detection Using YOLOv8

YOLOv8n is used as the detection backbone, pre-trained on the COCO dataset [2] and applied directly without any additional fine-tuning. Each detection produces a bounding box  $[x1, y1, x2, y2]$ , a confidence score, and a class ID. Detections are split into persons (class 0, confidence  $\geq 0.45$ ) and personal belongings (classes 24, 26, 28, 63, 67, 73 —confidence  $\geq 0.35$ ).

#### Person State Management

Each detected person is assigned a unique Person ID (PID). Persons are tracked across frames using centroid distance matching. If a new detection's centroid is within 120 pixels of an existing person's centroid, they are matched as the same person. Persons who do not appear in the current frame are immediately removed from the active list.

$$C = ((x1+x2)/2, (y1+y2)/2) \dots (1)$$

Match condition:  $d(C_{\text{new}}, C_{\text{existing}}) < 120 \text{ px} \dots (2)$

#### Object Tracking via IoU Matching

Personal belongings are tracked across frames using Intersection-over-Union (IoU) matching. If the IoU between a new detection and an existing tracked object exceeds 0.35, they are matched as the same object. Each object tracker stores its bounding box, class name, first-seen time, last-seen time, and abandoned flag.

$$\text{IoU}(A,B) = |A \cap B| / |A \cup B| \dots (3)$$

### Spatial Proximity Evaluation

For each tracked object, the system calculates the Euclidean distance between the object's centroid and the centroid of every active person. If at least one person is within 200 pixels, the object is considered attended. If no person is within this distance, the abandonment timer continues.

$$d(O,P) = \sqrt{(cx_O-cx_P)^2 + (cy_O-cy_P)^2} \dots \quad (4)$$

### Abandonment Classification

An object is classified as abandoned when two conditions are both satisfied simultaneously: (1) Object age >2.0 seconds, and (2) No person within 200 pixels of the object.

$$\text{Abandoned} = \text{True iff (Age} > 2.0s) \text{ AND (d\_min} \geq 200 \text{ px)} \dots \quad (5)$$

When the flag transitions from False to True, the event is recorded in the incident log with the object class, time, and bounding-box location.

### Alert Generation and Output

When an object is classified as abandoned, the system draws an orange bounding box and warning label showing the object class and how long it has been unattended. Objects in the monitoring state (1–2 s old) show a yellow box. All persons are shown with a green box labeled with their Person ID. The annotated video is saved as an MP4 file and a separate incident log records all abandonment events with timestamps and coordinates.

## IV. PROPOSED SYSTEM ARCHITECTURE

The proposed system is designed to automatically detect abandoned objects in real-time surveillance video without requiring any human involvement. The entire pipeline is built around a single deep learning model — YOLOv8 — which handles all object and person detection tasks. The output is then processed by five additional modules that together determine whether a detected object has been abandoned. The system processes video frame by frame, detects all persons and personal belongings, tracks them across frames, and checks whether any person is nearby.

### Object Detection Using YOLOv8

YOLOv8 is used as the detection backbone, pre-trained on the COCO dataset [2] and applied directly without any additional fine-tuning. Each detection produces a bounding box [x1, y1, x2, y2], a confidence score, and a class ID. Detections are split into persons (class 0, confidence  $\geq 0.45$ ) and personal belongings (classes 24, 26, 28, 63, 67, 73 — confidence  $\geq 0.35$ ).

### Person State Management

Each detected person is assigned a unique Person ID (PID). Persons are tracked across frames using centroid distance matching. If a new detection's centroid is within 120 pixels of an existing person's centroid, they are matched as the same person. Persons who do not appear in the current frame are immediately removed from the active list.

$$C = ((x1+x2)/2, (y1+y2)/2) \dots \quad (1)$$

$$\text{Match condition: } d(C_{\text{new}}, C_{\text{existing}}) < 120 \text{ px} \dots \quad (2)$$

### Object Tracking via IoU Matching

Personal belongings are tracked across frames using Intersection-over-Union (IoU) matching. If the IoU between a new detection and an existing tracked object exceeds 0.35, they are matched as the same object. Each object tracker stores its bounding box, class name, first-seen time, last-seen time, and abandoned flag.

$$\text{IoU}(A,B) = |A \cap B| / |A \cup B| \dots \quad (3)$$

### Spatial Proximity Evaluation

For each tracked object, the system calculates the Euclidean distance between the object's centroid and the centroid of every active person. If at least one person is within 200 pixels, the object is considered attended. If no person is within this distance, the abandonment timer continues.

$$d(O,P) = \sqrt{(cx_O-cx_P)^2 + (cy_O-cy_P)^2} \dots \quad (4)$$

### Abandonment Classification

An object is classified as abandoned when two conditions are both satisfied simultaneously: (1) Object age > 2.0 seconds, and (2) No person within 200 pixels of the object.

$$\text{Abandoned} = \text{True iff (Age} > 2.0s) \text{ AND (d\_min} \geq 200 \text{ px)} \dots \quad (5)$$

When the flag transitions from False to True, the event is recorded in the incident log with the object class, time, and bounding-box location.

#### Alert Generation and Output

When an object is classified as abandoned, the system draws an orange bounding box and warning label showing the object class and how long it has been unattended. Objects in the monitoring state (1–2 s old) show a yellow box. All persons are shown with a green box labeled with their Person ID. The annotated video is saved as an MP4 file and a separate incident log records all abandonment events with timestamps and coordinates.

### IV. METHODOLOGY

#### Data Collection and Input Setup

The input to the proposed system is a standard MP4 surveillance video file recorded from a fixed-angle CCTV camera in a public indoor environment such as a corridor, waiting area, or lobby. The system works with any standard resolution video, typically 720p (1280×720 pixels) or 1080p (1920×1080 pixels), recorded at 24 to 30 frames per second. No custom dataset was created — the system relies entirely on the pre-trained YOLOv8n model with standard COCO class labels.

#### Preprocessing and Frame Handling

Each frame is read in BGR format and passed directly to the YOLOv8 model without any additional resizing or normalisation. A frame counter is maintained throughout processing and the current timestamp in seconds is computed as:

$$\text{Timestamp (s)} = \text{Frame Number} / \text{Frames per Second} \dots (6)$$

This timestamp is used by the object tracker to calculate how long each detected object has been present in the scene. No data augmentation or background modelling is performed during preprocessing.

#### Object Detection Methodology

Each video frame is passed to the YOLOv8n model, which processes the entire frame in a single forward pass. The nano variant (YOLOv8n) is the smallest and fastest model in the YOLOv8 family, making it suitable for deployment on systems without a dedicated GPU. The raw output is a list of detection records, each containing six values: x1, y1, x2, y2,

confidence score, and class ID. Only detections that pass the confidence threshold are kept for further processing.

#### Tracking and Abandonment Logic

After detection filtering is complete for a given frame, the system updates both the person tracker and the object tracker independently. For person tracking, centroid coordinates are compared using Euclidean distance — if the closest existing centroid is within 120 pixels, the detection inherits the existing Person ID. For object tracking, bounding boxes are compared using IoU. If no match is found, a new ObjTracker instance is created recording the class name and first-seen timestamp.

$$\text{Abandoned} = \text{True iff (Age} > 2.0\text{s) AND (d\_min} \geq 200 \text{px)} \dots (7)$$

#### Output Generation Methodology

Once the abandonment decision has been made for all tracked objects in a frame, the system draws all visual annotations using OpenCV drawing functions. Annotations are rendered in order: person boxes first, then object monitoring boxes, then abandoned object alert boxes on top. Every annotated frame is written to the output video file using OpenCV's VideoWriter with the MP4V codec at the same resolution and frame rate as the input.

### V. EXPERIMENTAL RESULTS

#### Experimental Setup

Experiments were conducted using real surveillance-style video footage recorded in indoor public environments from a fixed overhead angle, simulating a standard CCTV installation. The system was executed entirely on a standard CPU-based machine to represent a realistic low-cost deployment scenario. No GPU acceleration was used during testing.

Component	Specification
Execution Environment	Google Colab (CPU runtime)
Programming Language	Python 3.10
Detection Model	YOLOv8n (COCO pretrained)
Video Processing	OpenCV 4.x
Input Resolution	1280 x 720 pixels (720p)
Frame Rate	24-30 FPS
Output Format	MP4 (MP4V codec)
Model Weights	yolov8n.pt (zero-shot)

Object Detection Results

The YOLOv8n model detected persons and personal belongings reliably across all test sequences. Persons were detected with confidence scores ranging from 0.48 to 0.91, well above the threshold of 0.45. Personal belongings were detected with confidence scores mostly between 0.38 and 0.82. The model correctly identified all six belonging categories present in the test videos.

Object Class	COCO ID	Detected	Avg Conf.
Person	0	Yes	0.72
Backpack	24	Yes	0.61
Handbag	26	Yes	0.58
Suitcase	28	Yes	0.55
Laptop	63	Yes	0.63
Mobile Phone	67	Partial	0.42
Book	73	Partial	0.39

Tracking Performance

The IoU-based object tracking mechanism maintained stable and consistent object identities throughout all test sequences. Abandoned objects being stationary, their bounding boxes in consecutive frames overlapped with IoU scores typically between 0.75 and 0.95. Person tracking using centroid distance matching also performed reliably with no identity switches observed during the main test sequences.

Scenario	Expected	System Result	OK?
Person drops bag, walks away (>2s)	Alert	Alert at 2.1s	Yes
Person sets bag down, stays nearby	No alert	No alert	Yes
Person steps away (<2s) then returns	No alert	Timer reset	Yes
Two bags left by two people	Two alerts	Two alerts	Yes
Person walks past abandoned bag	Alert clears	Alert cleared	Yes
Small object (phone) hidden	May miss	Not detected	Exp.

Abandonment Detection Results

The abandonment detection logic was tested across multiple scenarios, each designed to evaluate a specific aspect of the system's decision-making. The system correctly flagged objects as abandoned in all cases where a bag was left alone for more than 2 seconds with no person within 200 pixels.

Metric	Value	Notes
Person Detection Accuracy	91.3%	4 persons in frame
Object Detection Accuracy	87.6%	Lower for small objects
Abandonment Precision	88.9%	True alerts proportion
Abandonment Recall	85.2%	Real abandonments caught
False Positive Rate	11.1%	Brief owner departure
False Negative Rate	14.8%	Small/hidden objects
Avg Inference Time	~24 ms/frame	YOLOv8n on CPU
Processing Speed	~18 FPS	CPU only
Alert Response Delay	2.0-2.2 s	Abandonment to alert
Overall System Accuracy	89.4%	Detection + alert

System Performance Metrics

System performance was evaluated by measuring detection accuracy, alert precision, and processing speed across all test sequences. All values were measured on a standard CPU without GPU acceleration, representing a conservative deployment scenario.

State	Box Color	Label Format
Person detected	Green	Person N: Walking
Monitoring (1-2 s)	Yellow	Monitoring [Class]   X s
Abandoned (>2 s)	Orange	UNATTENDED [CLASS]   X s
Blinking alert banner	Orange text	!!! UNATTENDED OBJECT !!!

V. CONCLUSION

This paper presented an intelligent surveillance system that can automatically detect abandoned personal belongings in real-time video footage using the YOLOv8 deep learning model. By combining object detection, IoU-based tracking, and spatial proximity analysis into a single lightweight pipeline, the proposed system is able to identify unattended objects without any manual involvement.

The system was built using the YOLOv8n model in a zero-shot configuration, meaning no custom training data or fine-tuning was required. Experimental results showed an overall accuracy of 89.4%, with a person detection accuracy of 91.3% and an object detection accuracy of 87.6%. The system processed video at approximately 18 frames per second on a standard CPU, which is sufficient for near-real-time surveillance.

The system does have certain limitations, including occasional false alerts when an owner steps briefly outside the proximity radius, and difficulty detecting very small or partially hidden objects. These limitations will be addressed in future work. Overall, the proposed system demonstrates that an effective abandoned object detection solution can be built using publicly available tools and pre-trained models, without expensive hardware or custom datasets. It is ready for integration into existing CCTV infrastructure in airports, railway stations, shopping malls, and other public spaces where security monitoring is critical.

## REFERENCES

- [1] G. Jocher, A. Chaurasia, and J. Qiu, "Ultralytics YOLOv8," Ultralytics, 2023. [Online]. Available: <https://github.com/ultralytics/ultralytics>
- [2] T.-Y. Lin et al., "Microsoft COCO: Common objects in context," in Proc. ECCV, Zurich, Sep. 2014, pp. 740-755.
- [3] C. Stauffer and W. E. L. Grimson, "Adaptive background mixture models for real-time tracking," in Proc. IEEE CVPR, Jun. 1999, vol. 2, pp. 246-252.
- [4] M. Piccardi, "Background subtraction techniques: A review," in Proc. IEEE SMC, Oct. 2004, vol. 4, pp. 3099-3104.
- [5] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection," in Proc. IEEE CVPR, Jun. 2014, pp. 580-587.
- [6] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in Proc. NIPS, Dec. 2015, vol. 28, pp. 91-99.
- [7] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in Proc. IEEE CVPR, Jun. 2016, pp. 779-788.
- [8] J. Redmon and A. Farhadi, "YOLOv3: An incremental improvement," arXiv:1804.02767, Apr. 2018.
- [9] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, "Simple online and realtime tracking," in Proc. IEEE ICIP, Sep. 2016, pp. 3464-3468.
- [10] N. Wojke, A. Bewley, and D. Paulus, "Simple online and realtime tracking with a deep association metric," in Proc. IEEE ICIP, Sep. 2017, pp. 3645-3649.
- [11] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "YOLOv4: Optimal speed and accuracy of object detection," arXiv:2004.10934, Apr. 2020.
- [12] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, "YOLOX: Exceeding YOLO series detectors," arXiv:2107.08430, Jul. 2021.
- [13] A. Vaswani et al., "Attention is all you need," in Proc. NIPS, Dec. 2017, vol. 30, pp. 5998-6008.
- [14] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in Proc. IEEE CVPR, Jun. 2016, pp. 770-778.
- [15] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in Proc. ICLR, May 2015.
- [16] K. Lim et al., "CCTV violence detection using deep CNN and RNN," Electronics, vol. 10, no. 13, p. 1574, Jun. 2021.
- [17] G. Tripathi, K. Singh, and D. K. Vishwakarma, "CNNs for crowd behaviour analysis: A survey," The Visual Computer, vol. 35, no. 5, pp. 753-776, May 2019.
- [18] X. Zhou, D. Wang, and P. Krahenbuhl, "Objects as points," arXiv:1904.07850, Apr. 2019.
- [19] W. Liu et al., "SSD: Single shot multibox detector," in Proc. ECCV, Oct. 2016, pp. 21-37.
- [20] P. Dollar, R. Appel, S. Belongie, and P. Perona, "Fast feature pyramids for object detection," IEEE TPAMI, vol. 36, no. 8, pp. 1532-1545, Aug. 2014.
- [21] O. Barnich and M. Van Droogenbroeck, "ViBe: A universal background subtraction algorithm," IEEE Trans. Image Process., vol. 20, no. 6, pp. 1709-1724, Jun. 2011.
- [22] Z. Zivkovic and F. van der Heijden, "Efficient adaptive density estimation per image pixel," Pattern Recognit. Lett., vol. 27, no. 7, pp. 773-780, May 2006.
- [23] A. B. Chan, Z.-S. J. Liang, and N. Vasconcelos, "Privacy preserving crowd monitoring," in Proc. IEEE CVPR, Jun. 2008, pp. 1-7.
- [24] L. Leal-Taixe et al., "MOTChallenge 2015: Towards a benchmark for multi-target tracking," arXiv:1504.01942, Apr. 2015.
- [25] Y. Zhang et al., "ByteTrack: Multi-object tracking by associating every detection box," in Proc. ECCV, Oct. 2022, pp. 1-21.
- [26] A. Yilmaz, O. Javed, and M. Shah, "Object tracking: A survey," ACM Comput. Surv., vol. 38, no. 4, Dec. 2006.
- [27] J. Candamo et al., "Understanding transit scenes: A survey on human behavior-recognition algorithms," IEEE Trans. ITS, vol. 11, no. 1, pp. 206-224, Mar. 2010.
- [28] T. Ko, "A survey on behavior analysis in video surveillance for homeland security applications," in Proc. IEEE AIPR, Oct. 2008, pp. 1-8.
- [29] S. Patil and S. Kulkarni, "Deep learning-based abandoned object detection in surveillance videos: A survey," Multimedia Tools and Applications, vol. 81, no. 14, pp. 19743-19769, Jun. 2022.