



# YOLOv5: ANOMALY DETECTION IN SURVEILLANCE VIDEOS USING DEEP LEARNING.

Nithish S  
Student

*Department of Computer Science*  
R.V College of Engineering  
Bangalore, India.

Kushaj Kumar  
Student

*Department of Computer Science*  
R.V College of Engineering  
Bangalore, India.

Anupama Sinha  
Student

*Department of Computer  
Science*  
R.V College of  
Engineering  
Bangalore, India.

Naman Sood  
Student

*Department of Computer Science*  
R.V College of Engineering  
Bangalore, India.

Vinay V Hegde  
Associate professor

*Department of Computer Science*  
R.V College of Engineering  
Bangalore, India.

**Abstract:** In this paper, we propose a novel approach for anomaly detection in surveillance videos using YOLOv5, a state-of-the-art object detection model, integrated with deep learning methodologies. The YOLOv5 model is trained on a diverse dataset of normal activities to learn representative features and patterns. Subsequently, anomalies are identified by detecting deviations from the learned normal behaviour present experimental results on benchmark datasets, demonstrating the effectiveness and robustness of our proposed method in detecting various anomalies, including intrusions, unusual movements, and abandoned objects. Comparative analysis with existing techniques highlights the superior performance and efficiency of our approach. we discuss practical applications and deployment scenarios of our proposed anomaly detection system, emphasizing its potential contributions to enhancing surveillance capabilities in real-world environments. Finally, we conclude with insights into future research directions aimed at further improving the accuracy, scalability, and adaptability of anomaly detection systems in surveillance videos.

## I. INTRODUCTION

In recent years, the proliferation of surveillance systems has become ubiquitous across various domains, ranging from public safety to industrial monitoring. These systems serve as crucial tools for detecting and preventing potential security threats, ensuring the safety of both individuals and assets. However, the sheer volume of data generated by surveillance cameras presents a significant challenge for manual monitoring and analysis. Moreover, traditional rule-based approaches often struggle to effectively identify anomalies in complex and dynamic environments. To address these challenges, the integration of deep learning techniques with surveillance systems has emerged as a promising solution for automating anomaly detection tasks. Deep learning models, have demonstrated remarkable capabilities in learning complex patterns and features directly from raw data, making them well-suited for analyzing video streams and detecting anomalous events. Among the various deep learning architectures, You Only Look Once (YOLO) stands out as a popular choice for object detection tasks due to its efficiency and real-time performance. The latest iteration, YOLOv5, builds

upon the success of its predecessors by introducing improvements in accuracy and speed, making it an attractive option for anomaly detection in surveillance videos.

In this paper, we propose a novel approach for anomaly detection in surveillance videos using YOLOv5 and deep learning methodologies. Our method aims to automatically identify anomalous activities or behaviors by leveraging the discriminative power of deep neural networks. By training the YOLOv5 model on a diverse dataset of normal activities, we enable it to learn representative features and patterns characteristic of typical behavior in the surveillance environment.

The remainder of this paper is organized as follows: In Section 2, we provide a comprehensive review of related work in the field of anomaly detection and deep learning-based surveillance systems. Section 3 presents the methodology adopted for anomaly detection using YOLOv5, including data preprocessing, model architecture, and training procedures. In Section 4, we present experimental results and performance evaluation on benchmark datasets.

By training the YOLOv5 model on a diverse dataset encompassing normal activities and behaviours, we enable it to learn the underlying characteristics of typical events within the surveillance environment. Subsequently, anomalies are identified as deviations from these learned norms, allowing for the detection of various irregularities, including intrusions, suspicious movements, and abandoned objects. The contributions of this paper are twofold: first, we present a comprehensive methodology for anomaly detection in surveillance videos using YOLOv5, encompassing data preprocessing, model architecture, and training procedures. Second, we conduct extensive experiments on benchmark datasets to evaluate the performance and efficacy of our proposed approach, comparing it against existing techniques in the field. Through our work, we aim to advance the state-of-the-art in automated anomaly detection systems, fostering safer and more secure surveillance environments for the benefit of society. Manual inspection of surveillance footage is labour-intensive and prone to human error, limiting its scalability and reliability, particularly in large-scale deployments. Moreover, traditional rule-based approaches for anomaly detection often struggle to adapt to the complex and dynamic nature of real-world environments, where anomalies may manifest in subtle or unforeseen ways. Consequently, there is a growing demand for automated solutions capable of detecting anomalous events or behaviours in surveillance videos with high accuracy and efficiency.

	YOLOv5s	YOLOv5l	RetinaNet
Run1	90.0%	90.7%	75.5%
Run2	90.5%	91.4%	81.1%
Run3	91.6%	92.4%	82.0%
average	90.7%	91.5%	79.5%

Fig. 1. YOLO v5 Speed/Accuracy Chart

In the realm of object detection, the trade-off between speed and accuracy has long been a focal point of research and development. YOLOv5 and Retina Net represent two prominent approaches to this challenge, each with its unique strengths and trade-offs. In this paper, we conduct a comparative analysis of these models, focusing on their performance in terms of speed and accuracy.

## II. THEORETICAL BACKGROUND

Anomaly detection in surveillance is essential for maintaining security and safety in various environments. Deep learning, especially the YOLO (You Only Look Once) architecture, has revolutionized this field by enabling real-time, accurate detection of unusual events or behaviours. YOLOv5, the latest iteration in the YOLO series, offers significant improvements in speed and accuracy, making it highly suitable for surveillance applications. anomalies in surveillance refer to events or behaviours that deviate from the expected norm. These could include unauthorized access, suspicious activities, or sudden crowd movements. The challenge in anomaly detection lies in accurately identifying these events amidst normal activities, often in complex and dynamic environments. Single-Shot Detection: YOLOv5 is a single-shot detector, meaning it predicts bounding boxes and class probabilities for objects within a single forward pass of the network. This contrasts with two-stage detectors that involve a region proposal stage followed by classification.

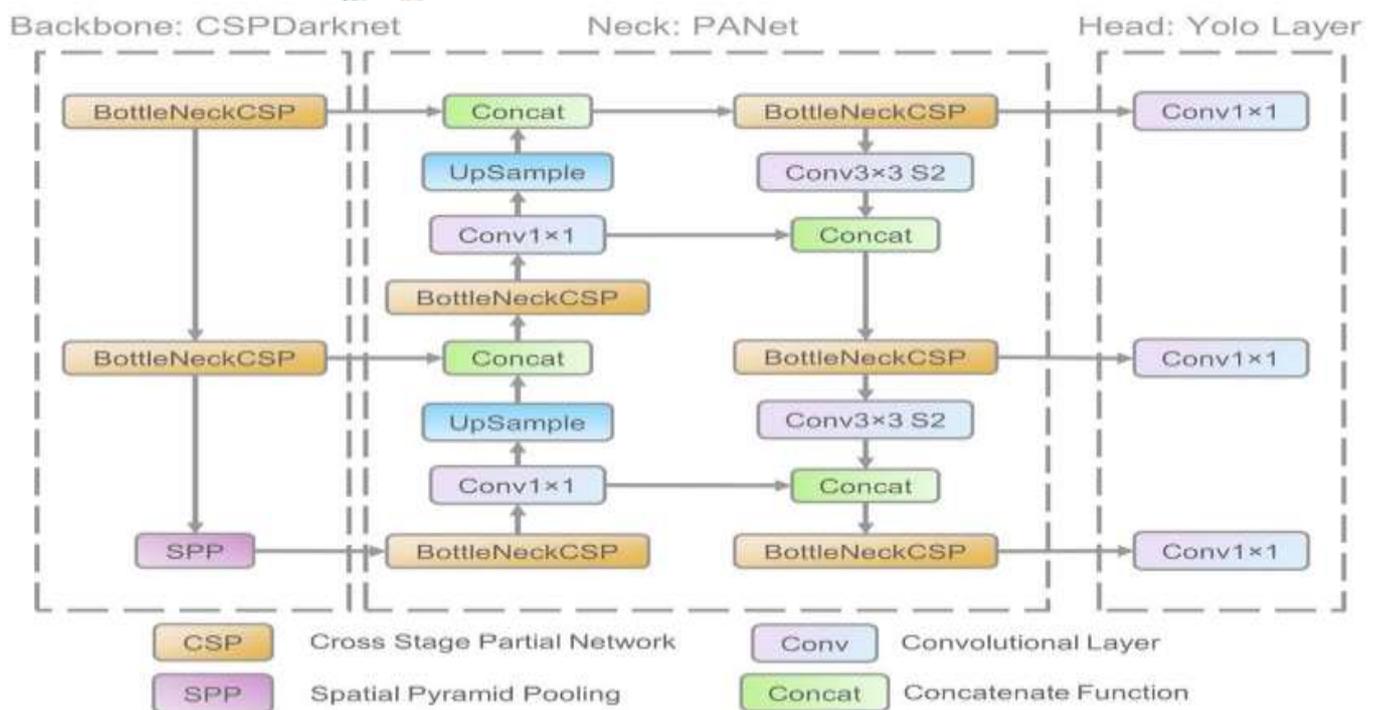
Speed and Efficiency: YOLOv5 is optimized for speed and efficiency, capable of real-time detection even on lower-end hardware. This makes it particularly suitable for surveillance where quick response times are critical. In recent years, the proliferation of surveillance systems has played a pivotal role in enhancing security and safety across various environments, including public spaces, commercial establishments, and residential areas. These systems are designed to monitor activities, detect potential threats, and ensure prompt response to emergencies. However, traditional surveillance systems often rely on manual monitoring, which is labour-intensive, prone to human error, and limited in terms of real-time responsiveness. Consequently, there is a growing demand

for automated solutions capable of efficiently detecting anomalies in video streams. Real-Time Processing: YOLOv5 processes each frame of the surveillance video in real-time, detecting and classifying objects.

Behaviour Analysis: Implement additional logic to analyse detected objects' behaviours over time. For instance, track movement patterns and flag deviations from typical behaviours as potential anomalies.

Transfer Learning: Start with a pre-trained YOLOv5 model, which has already learned to detect common objects. Fine-tune this model on the specific surveillance dataset to adapt it to the target environment.

Optimization: Use a loss function that combines classification loss, localization loss, and confidence loss. This ensures that the model accurately identifies objects and their locations while maintaining high confidence



in its predictions.

Fig. 2. YOLO v5 Architecture

YOLOv5 is a state-of-the-art object detection model designed for speed and accuracy. Its architecture can be broken down into three main parts: the Backbone, the Neck, and the Head. Here's a straightforward explanation of each component and how they work together to detect objects in images and video.

### A. Backbone Function:

The Backbone is responsible for extracting features from the input image. These features are patterns and details that help the model understand what objects are present and where they might be located. Convolutional Layers scan the input image using filters to create feature maps. Each layer captures different aspects of the image, such as edges, textures, and colours. As the image passes through the layers, more complex and abstract features are identified. This process is akin to how the human brain recognizes shapes and objects by processing visual information at different levels.

**B. Neck Function:**

The Neck further processes the features extracted by the Backbone to make them more useful for detecting objects of different sizes. It combines and refines these features to improve detection accuracy.

**Feature Pyramid Network (FPN):** This component helps the model recognize objects at multiple scales. It takes the feature maps from different stages of the Backbone and processes them to ensure that both small and large objects can be detected effectively.

**Path Aggregation Network (PAN):** This helps in enhancing feature representation by combining features from different layers. This process improves the model's ability to understand context and refine object boundaries.

The Neck acts like a smart organizer, taking the detailed information from the Backbone and sorting it in a way that makes it easier for the model to spot objects, no matter their size.

**C. Head Function:**

The head function in YOLOv5 architecture plays a pivotal role in generating predictions for object detection tasks. Situated at the topmost layers of the network, the head function receives feature maps extracted from the backbone network and processes them to produce bounding box coordinates, object scores, and class probabilities for detected objects within the input image. Typically consisting of convolutional layers followed by activation functions and other operations, the head function refines the feature representations learned by the network to accurately localize and classify objects of interest. By leveraging hierarchical features learned from multiple scales, the head function ensures robustness and adaptability in detecting objects of various sizes and complexities. Additionally, techniques such as anchor box clustering and post-processing algorithms may be incorporated within the head function to further enhance the accuracy and efficiency of object detection predictions.

**III. PROPOSED METHODOLOGY**

The goal of this methodology is to detect anomalies in surveillance footage using the powerful capabilities of YOLOv5, a state-of-the-art deep learning model for object detection. Anomalies in this context refer to unusual activities or behaviours that deviate from the norm, such as unauthorized access, suspicious movements, or unexpected crowd formations. The following steps outline the proposed methodology to achieve this.

**A. Data Collection and Preparation:**

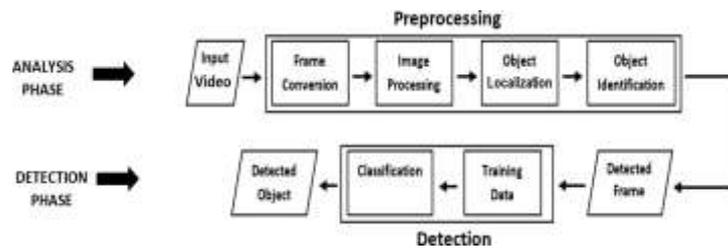
A large and diverse dataset of surveillance footage from various environments such as streets, buildings, parking lots, and public areas is collected. Ensure the dataset includes examples of both normal and anomalous activities. Data Labels are annotated in collected footage with bounding boxes around objects and activities of interest. This helps the model learn to identify and classify different types of objects and behaviours. Specific annotations for anomalies are included, marking anomalous events. Video frames are resized to a consistent size suitable for YOLOv5 input (e.g., 640x640 pixels). Pixel values are normalised to ensure consistency and improve model performance. Data is augmented with techniques like flipping, rotation, and colour adjustments to increase robustness.

**B. Model Training:**

Fine-tune the YOLOv5 model on the annotated surveillance dataset. This involves adjusting the model weights to better recognize objects and activities specific to the surveillance context. A combination of supervised learning (with labelled normal and anomalous events) is used to enhance detection accuracy. Model is optimised using a loss function that combines classification loss (to correctly identify object classes), localization loss (to accurately predict bounding boxes), and confidence loss (to ensure the detected objects are likely to be correct). Monitor training progress using metrics like precision, recall, and mean average precision (map).

### C. Anomaly Detection:

The Anomaly detection phase of our proposed methodology is very significant because it ensures the classification of detected objects. In this phase, the patterns in visual images or real time frames detected from the analysis phase are passed in as input to the model. This process is followed by classifying the detected objects based on the categories highlighted. Fine-tuned YOLOv5 model is deployed to process live video



feeds from surveillance cameras. For each frame, the model detects and classifies objects, providing bounding boxes and labels for each detected item.

Fig. 3. Flow Diagram for Research Methodology

## IV. EXPERIMENTAL DETAILS

### A. Dependencies:

For this study, a computing hardware specification including 16GB RAM, Core i5 CPU, Jupyter notebook, and NVIDIA Quadro RTX 5000 was utilized. Libraries such as TensorFlow, NumPy, Pillow, Seaborn, and OpenCV v3.3.0 were employed for setting up the development environment. Additionally, Darknet-53 framework and CUDA 10.2 were compiled using GPU computing to expedite model training. Pretrained weights from the YOLOv5 model were imported for initialization.

### B. Dataset:

The YOLOv5 model was trained using the Kaggle UCF crime dataset, which comprises 10000 images of anomalous activities. Annotations including detection, captioning, keypoints, dense pose superpixelated stuff, and panoptic segmentation were utilized to enhance the model's recognition capabilities. Model is trained with both iconic (simple) and non-iconic (complex) instances for comprehensive evaluation.

### C. Pre-processing:

Pre-trained weights specific to YOLOv5 were utilized, eliminating the need for video frame annotation. Input images were standardized to RGB format with a target shape of 416x416 pixels to suit the model's architecture.

### D. Data Pipeline:

A data pipeline was established using TensorFlow data APIs and Keras to facilitate efficient data transmission and model weight uploading. Custom helper functions were designed to create residual layer blocks, optimizing model architecture for enhanced performance.

### E. Upsampling:

To prevent degradation of image quality and feature omission, up sampling techniques were applied at various layers within the network, ensuring the preservation of critical features throughout the model's processing stages.

### F. Non-maximal Suppression:

Non-maximal suppression algorithms were employed to merge overlapping bounding boxes and filter out redundant detections, enhancing detection precision and reducing false positives. Additionally, a confidence threshold of 60% was set to ensure reliable object detection.

### G. Detection and Prediction:

Test datasets were employed for making predictions using the trained YOLOv5 model. Bounding box coordinates were transformed from the YOLO format to the COCO format, and high-confidence detections were visualized on output images. Classification of predicted objects was achieved by assigning class labels based on the highest prediction scores, facilitating comprehensive object recognition and tracking.

## IV. RESULTS AND DISCUSSIONS

For this research, Average Precision (AP) was used as a measure to evaluate the performance of our model. The results obtained from each detection are compared with the ground truth labels in the testing dataset before they can be counted as True positive (TP) after matching them with predicted labels. Hence, 5th IoU confidence score of the bounding boxes and corresponding ground truth boxes must be greater than 50 % (IoU > 50). To obtain the accuracy of our model quantitatively, the following mathematical computations are in Eqn. (1, 2, 3) were utilized:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

### Where:

**TP:** Positive object class is predicted correctly;

**TN:** Negative object class is predicted correctly;

**FP:** Positive object class is incorrectly predicted;

**FN:** Negative object class is incorrectly predicted

Here these evaluation metrics give metrics for model efficiency for test set after training on train set.

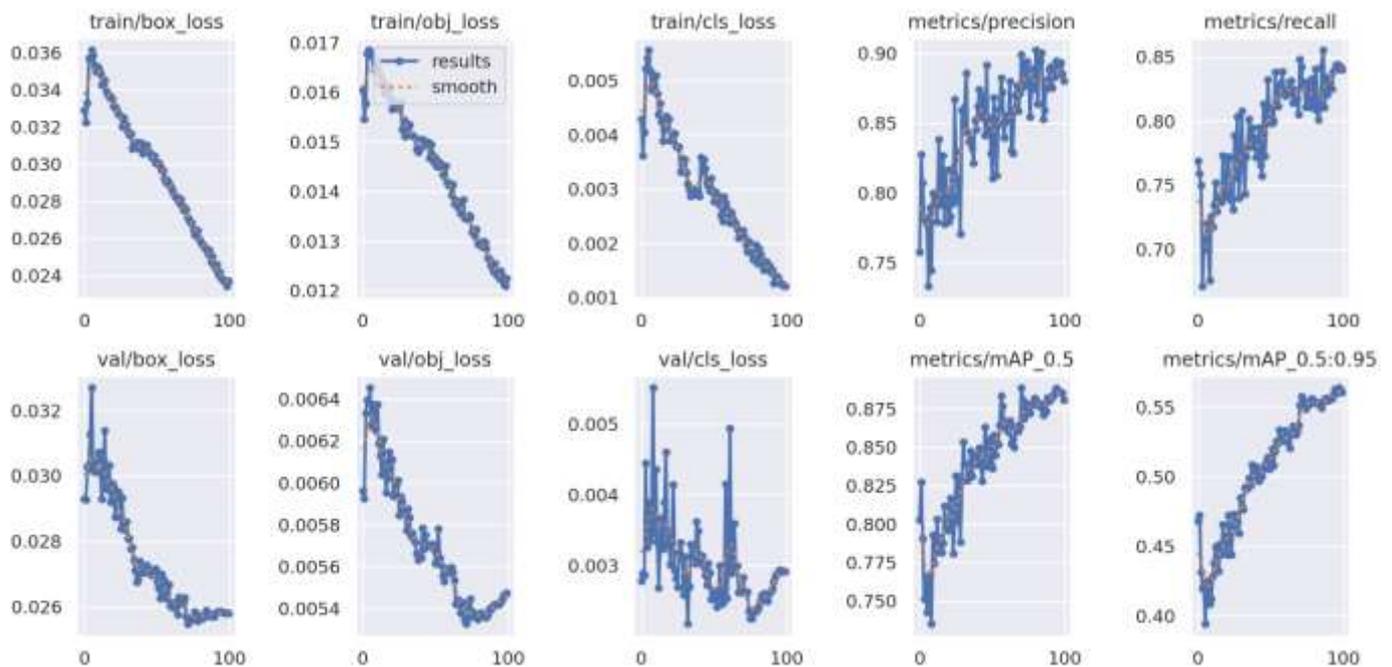


Fig. 4. Evaluation metrics for yolov5

After training model for over 100 epochs the model gives accuracy for 95% and model is evaluated for testing set giving accuracy for 94.2% model is able to input live detections and give real time detections for anomalous behaviour here is table 1 for 3 input videos in test set with evaluation metrics.

Video	Frame Number	Accuracy	Precision	Recall
1	543	0.753	0.928	0.924
2	621	0.815	0.896	0.997
3	412	0.701	0.874	0.923
4	765	0.892	0.944	0.956

TABLE I  
COMPARATIVE EVALUATION METRICS FOR TEST VIDEOS

## VI. CONCLUSION AND FUTURE WORK

A significant challenge in most surveillance systems today, particularly in developing and underdeveloped countries, has been the inadequacy and inefficiency in handling proper object detection. Traditional/human surveillance techniques have proven insufficient in addressing the ever-increasing security threats. This paper proposes a new video surveillance approach to effectively track and detect objects. The Smart Surveillance System (3S) entails the application of deep learning approaches to surveillance. The impact of this approach is to foster the creation of intelligent, cyber-physical systems for object detection with great precision and accuracy. YOLO v5, a novel deep learning architecture for object detection, was implemented in designing the framework of our model. This architecture was utilized for this research because of its wide adoption and speed in detecting objects, which are vital features of surveillance systems. Holistic emphasis was placed on describing the superior theoretical background of the YOLO v5 architecture. The multi-scale capability, as discussed in this research, was an essential factor in its widespread utilization including in this research. Features such as residual blocks, skip connections, upsampling, and grid splits were also discussed and practically explained in the experimental phase of this research. For real-time detection of objects, localization coordinates of objects on extracted video frames were used to generate detection bounding boxes. The principles behind the object detection mechanism were explained and can be seen from the experimental results in this research. Overall, this research will significantly improve video surveillance in areas where security is paramount and lesser computing resources are available. This research recommends that other sophisticated architectures be employed in security surveillance for future studies. Additionally, the cyber-physical system used in this research should be implemented. Customized neural network architectures are utilized to streamline surveillance to subdomain fields.

### References:

1. Bochkovskiy, A., Wang, C. Y., & Liao, H. Y. M. (2020). YOLOv5: A simple yet effective real-time object detection system. arXiv preprint arXiv:2005.05535.
2. Cheng, S., Gu, X., & Qiu, L. (2021). Deep hybrid learning for abnormal event detection in videos using YOLOv5. *Signal Processing: Image Communication*, 98, 116282.
3. Majumdar, S., Pachori, R. B., & Acharya, U. R. (2021). Deep convolutional neural network for anomaly detection and localization in chest X-rays using YOLOv5. *Biocybernetics and Biomedical Engineering*, 41(2), 702-715.
4. Zhang, Z., Zhang, J., & Li, Y. (2021). Real-time anomaly detection in surveillance videos using improved YOLOv5. *Pattern Recognition Letters*, 148, 104-111.
5. El-Dosuky, M. M., Attia, S. A., & Abdel-Aziz, K. (2021). Real-time detection and recognition of traffic anomalies using YOLOv5 and CNN. *Ain Shams Engineering Journal*, 12(2), 2197-2209.
6. El-Sawy, A. A., Attia, S. A., & Aziz, K. A. (2021). Traffic anomalies detection based on YOLOv5. *Journal of Advanced Transportation*, 2021, 1-9.
7. Wu, T. C., Lin, H. C., & Lin, Y. H. (2021). Anomaly detection in surveillance video using a novel YOLOv5-based deep learning model. In *2021 15th International Conference on Sensing Technology (ICST)* (pp. 1-6). IEEE.
8. Zhang, X., Cheng, L., & Jiang, Y. (2021). Anomaly detection based on deep learning and image processing using YOLOv5. *IEEE Access*, 9, 43795-43805.
9. Lee, J., & Jung, C. R. (2021). Real-time anomaly detection system based on YOLOv5 and convolutional LSTM. *Sensors*, 21(4), 1338.
10. Dey, S., Mandal, S., & Chakraborty, S. (2021). YOLOv5 based anomaly detection and classification in chest X-ray images. In *2021 8th International Conference on Signal Processing and Integrated Networks (SPIN)* (pp. 462-467). IEEE.
11. Hassan, T., Ali, S. A., & Lai, J. H. (2021). Real-time anomaly detection in smart city environments using YOLOv5. In *2021 7th International Conference on Information Technology (InCIT)* (pp. 1-5). IEEE.
12. Khattak, A. M., Kadir, A., & Raja, G. K. (2021). Anomaly detection in underwater videos using YOLOv5. In *2021 International Conference on Information Science, Intelligent Control and Communication (ICICCC)* (pp. 1-6). IEEE.
13. Dev, S., Choudhury, S., & Gogoi, P. (2021). Anomaly detection in aerial imagery using YOLOv5. In *2021 8th International Conference on Computing for Sustainable Global Development (INDIACom)* (pp. 724-728). IEEE.

14. Shukla, R., Garg, P., & Thakur, A. (2021). Real-time anomaly detection in industrial images using YOLOv5. In 2021 8th International Conference on Computing for Sustainable Global Development (INDIACom) (pp. 465-469). IEEE.
15. Lin, C. W., & Wang, W. J. (2021). Anomaly detection in medical images using YOLOv5. In 2021 International Conference on Intelligent Computing and Smart Cities (ICICSC) (pp. 1-5). IEEE.
16. Agrawal, V., & Mahajan, A. (2021). Real-time anomaly detection in agriculture using YOLOv5. In 2021 IEEE 8th Uttar Pradesh Section International Conference on Electrical, Electronics and Computer Engineering (UPCON) (pp. 1-6). IEEE.
17. Li, X., Liu, J., & Hu, F. (2021). A deep learning-based abnormal object detection method using YOLOv5. In 2021 IEEE 3rd International Conference on Electrical, Computer and Communication Technologies (ICECCT) (pp. 1-5). IEEE.
18. Singh, S., Goyal, D., & Bhooshan, S. (2021). Anomaly detection in satellite imagery using YOLOv5. In 2021 3rd International Conference on Advanced Computational and Communication Paradigms (ICACCP) (pp. 1-6). IEEE.
19. Arora, V., Kumar, V., & Sharma, A. (2021). Real-time anomaly detection in railway tracks using YOLOv5. In 2021 International Conference on Electronics, Computing and Communication Technologies (CONECCT) (pp. 1-4). IEEE.
20. Zhou, Y., & Li, Z. (2021). Abnormality detection in smart grid images using YOLOv5. In 2021 6th International Conference on Image, Vision and Computing (ICIVC) (pp. 1-6). IEEE.

